

## LA-UR-19-21274

Approved for public release; distribution is unlimited.

Title: Accelerated modeling of atomistic physics with machine learning

Author(s): Smith, Justin Steven; Lubbers, Nicholas Edward; Barros, Kipton Marcos; Nebgen, Benjamin Tyler; Tretiak, Sergei; Germann, Timothy Clark; Fensin, Saryu Jindal; Roitberg, Adrian E.; Isayev, Olexandr; zubatyuk, roman; Burakovsky, Leonid; Devereux, Christian; Ranashingha, Kavindri; Suwa, Hidemaro; Batista, Christian; Chern, Gai-Wei

Intended for: Machine Learning for Computational Fluid and Solid Dynamics, 2019-02-19/2019-02-21 (Santa Fe, New Mexico, United States)

Issued: 2019-02-19

---

**Disclaimer:**

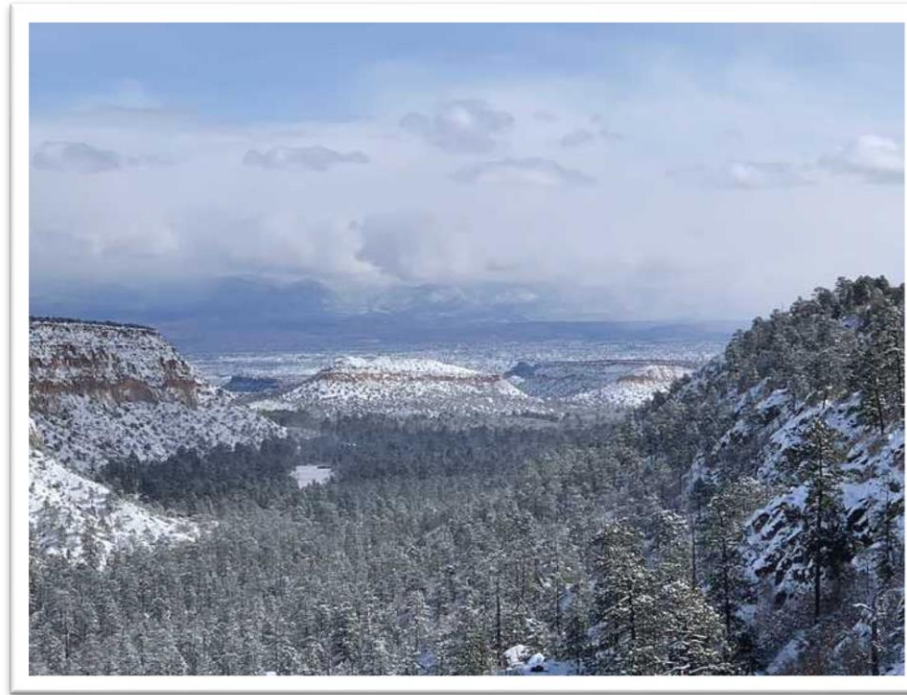
Los Alamos National Laboratory, an affirmative action/equal opportunity employer, is operated by Triad National Security, LLC for the National Nuclear Security Administration of U.S. Department of Energy under contract 89233218CNA000001. By approving this article, the publisher recognizes that the U.S. Government retains nonexclusive, royalty-free license to publish or reproduce the published form of this contribution, or to allow others to do so, for U.S. Government purposes. Los Alamos National Laboratory requests that the publisher identify this article as work performed under the auspices of the U.S. Department of Energy. Los Alamos National Laboratory strongly supports academic freedom and a researcher's right to publish; as an institution, however, the Laboratory does not endorse the viewpoint of a publication or guarantee its technical correctness.

# Accelerated modeling of atomistic physics with machine learning

## LANL Team

Benjamin Nebgen  
Kipton Barros  
Saryu Fensin  
Tim Germann  
Leonid Burakovsky  
Nicholas Lubbers  
Sergei Tretiak

Justin S. Smith



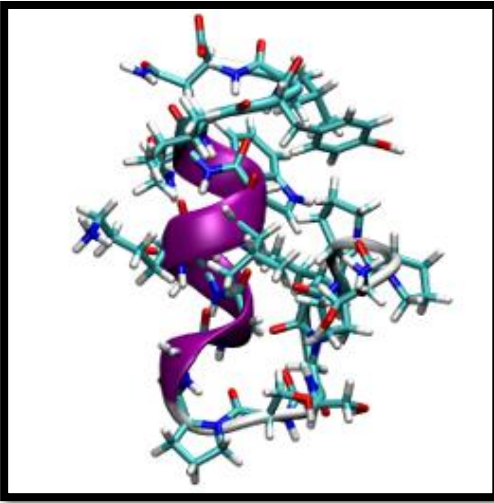
## Collaborators

Olexandr Isayev  
Adrian Roitberg  
Roman Zubatyuk  
Christian Devereux  
Kavindri Ranashingha  
Hidemaro Suwa  
Christian Batista  
Gia-Wei Chern

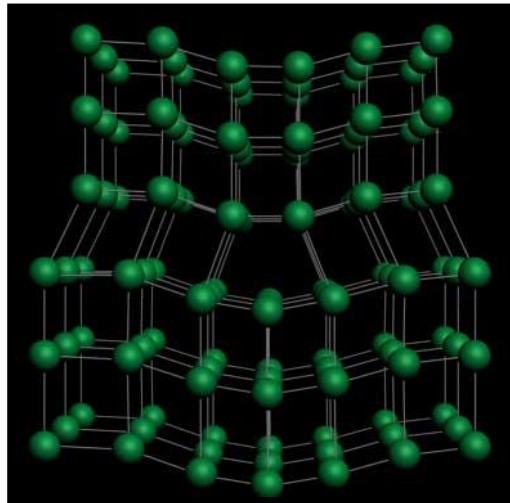
**Machine Learning for Computational Fluid and Solid Dynamics**

# Molecular (atomistic) dynamics

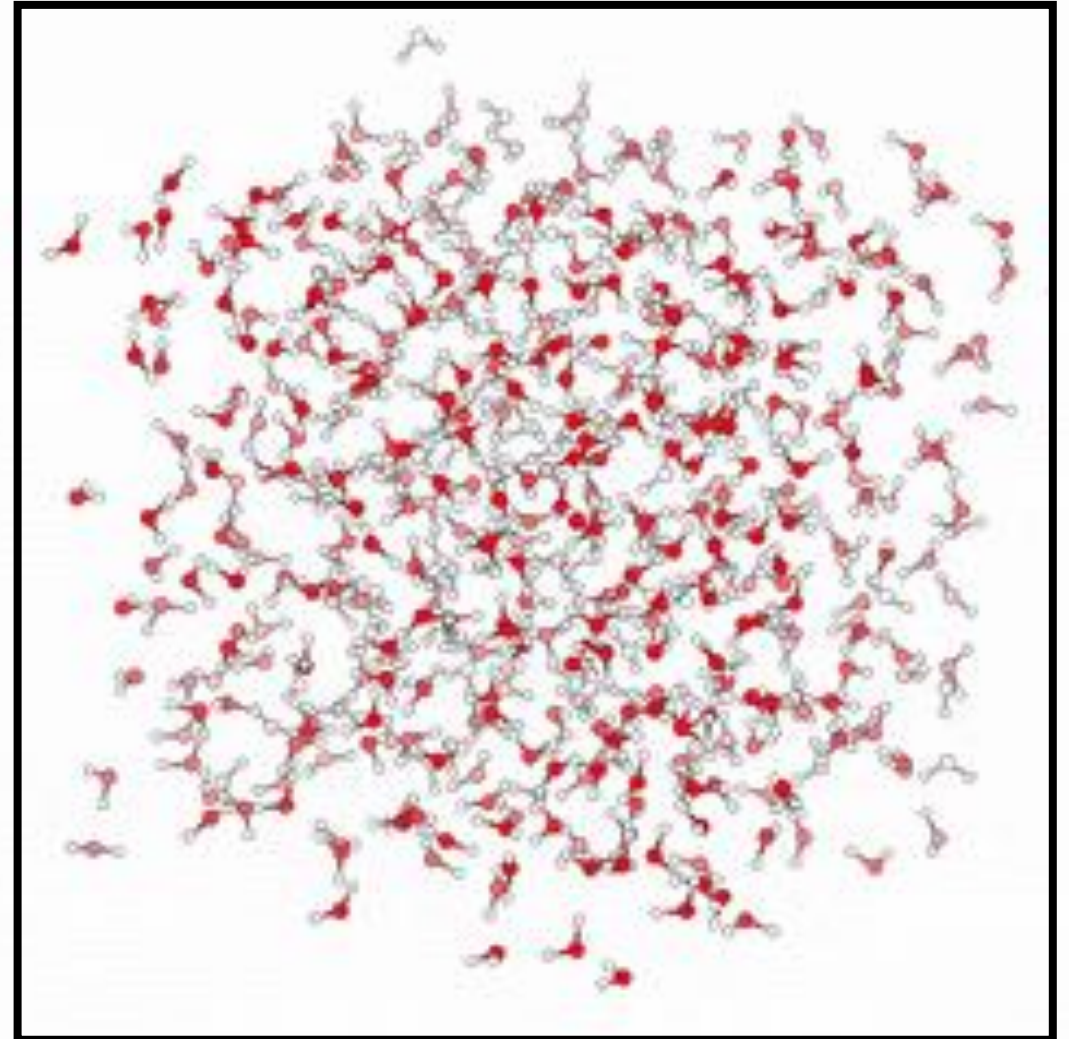
Proteins



Materials



Liquids



Energy

$$E[\mathbf{r}_1, \mathbf{r}_2, \dots]$$

Force

$$\mathbf{f}_i = -\nabla_i E$$

Dynamics

$$m \frac{d^2 \mathbf{r}_i}{dt^2} = \mathbf{f}_i$$

In principle, requires  
a quantum  
mechanical  
calculation at *each*  
time step!

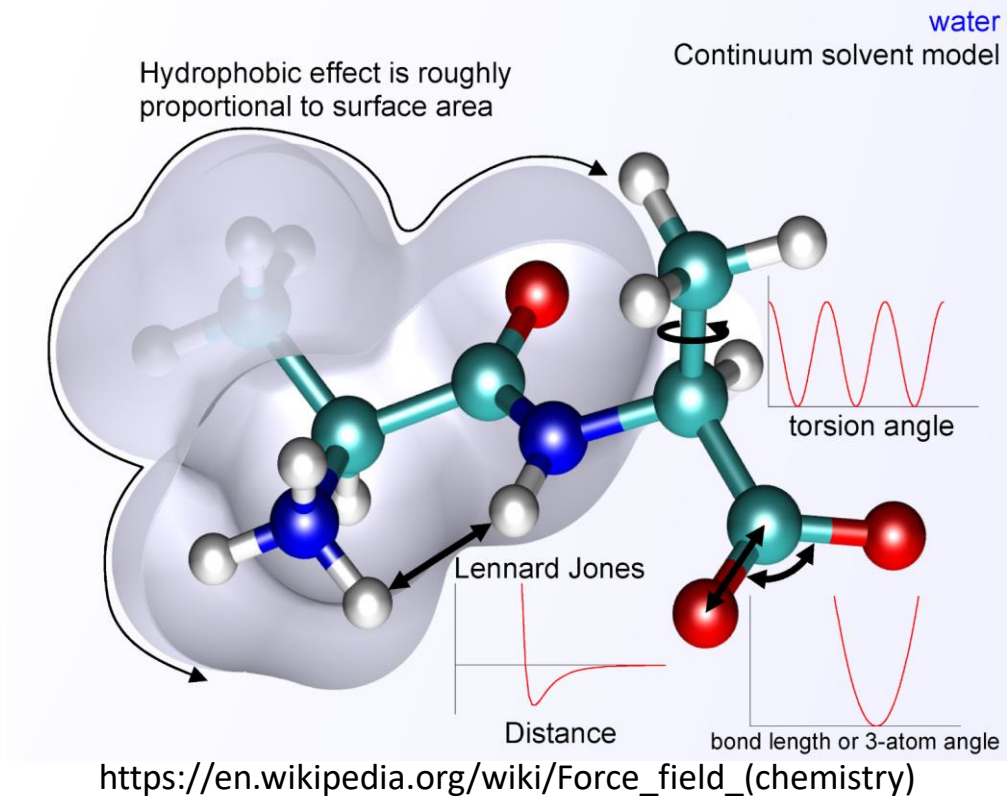
# Molecular mechanics/Classical force field

Pros:

- Computationally efficient
- Accurate on systems in fitting set

Cons:

- Not very transferable
- Non-reactive
- Difficult reparameterization



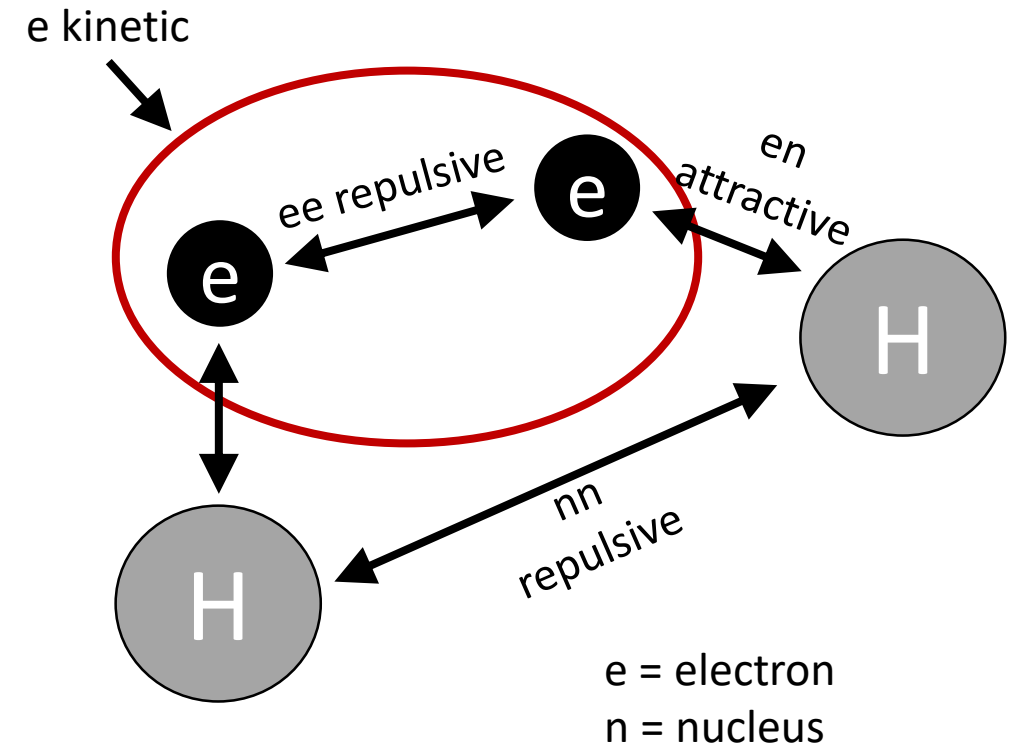
# The Electronic Schrödinger Equation (QM)

Pros:

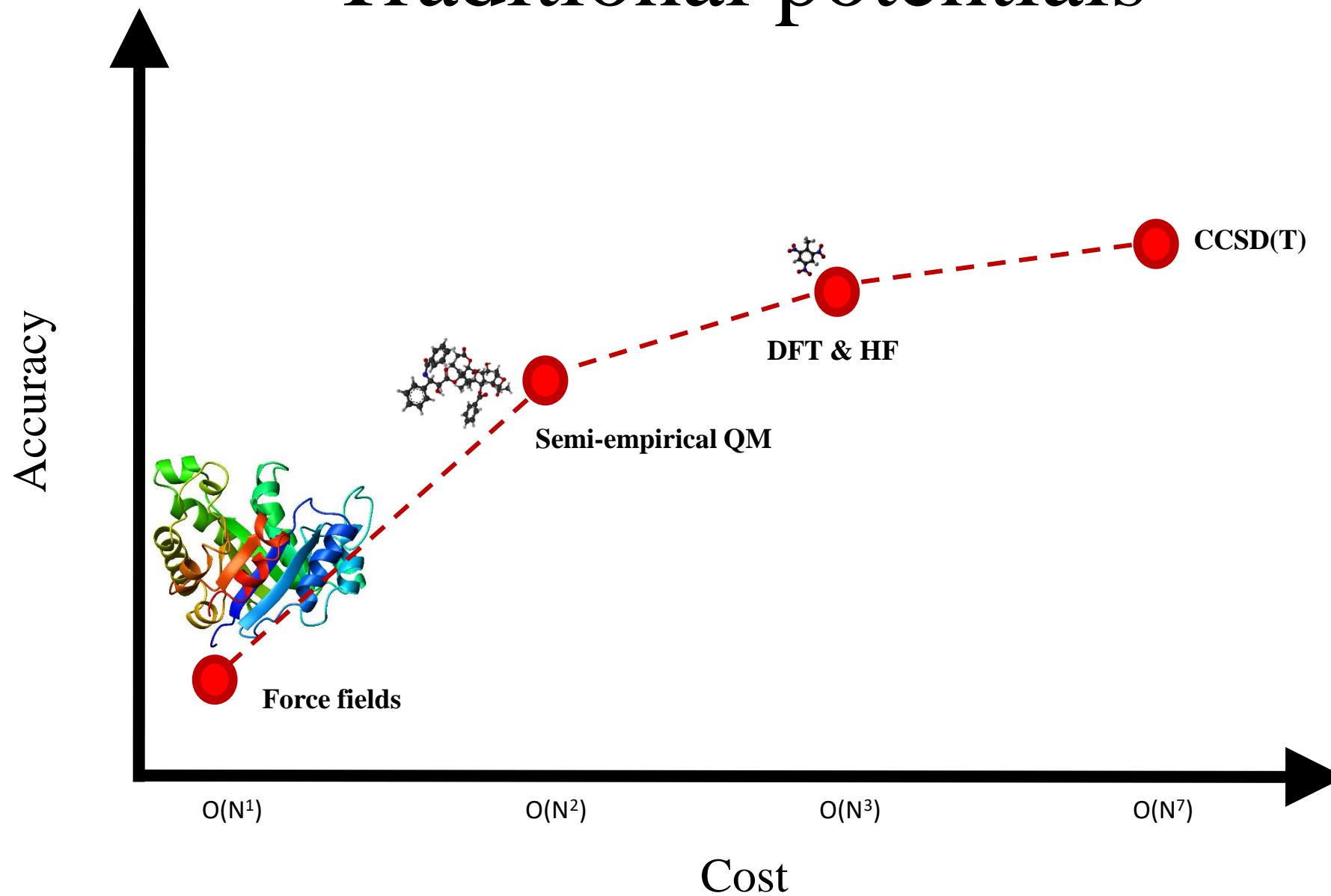
- Transferable
- Accurate

Cons:

- Computationally demanding



# Traditional potentials





# A **potential** solution

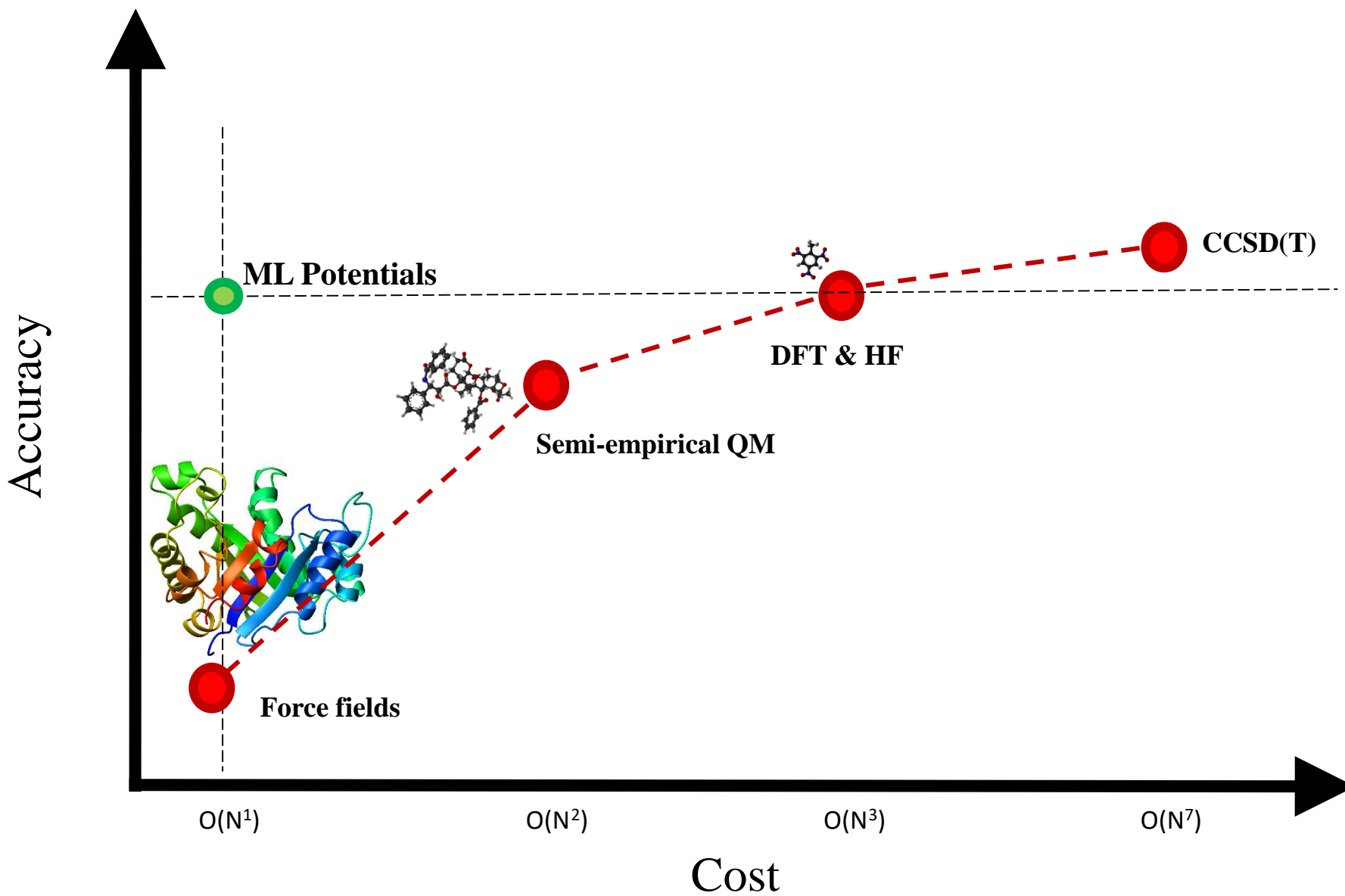
**Solution:** Develop an empirical potential that is accurate, fast and parametrizes itself

- ❑ *Machine learning* provides methods that fit this need
- ❑ Prior neural network potentials\* (NNP) for organic molecules and materials...
  - are trained to specific molecules or phases of a material
  - are non-transferable
- ❑ Our goal: build general and accurate ML potentials

Combine big data and deep learning concepts and a new molecular representation to produce *accurate, transferable, and extensible* NNPs.

\*Behler, J.; *Angew. Chemie Int. Ed.* **2017**, 56 (42), 12828–12840.

# Where does ML fit?

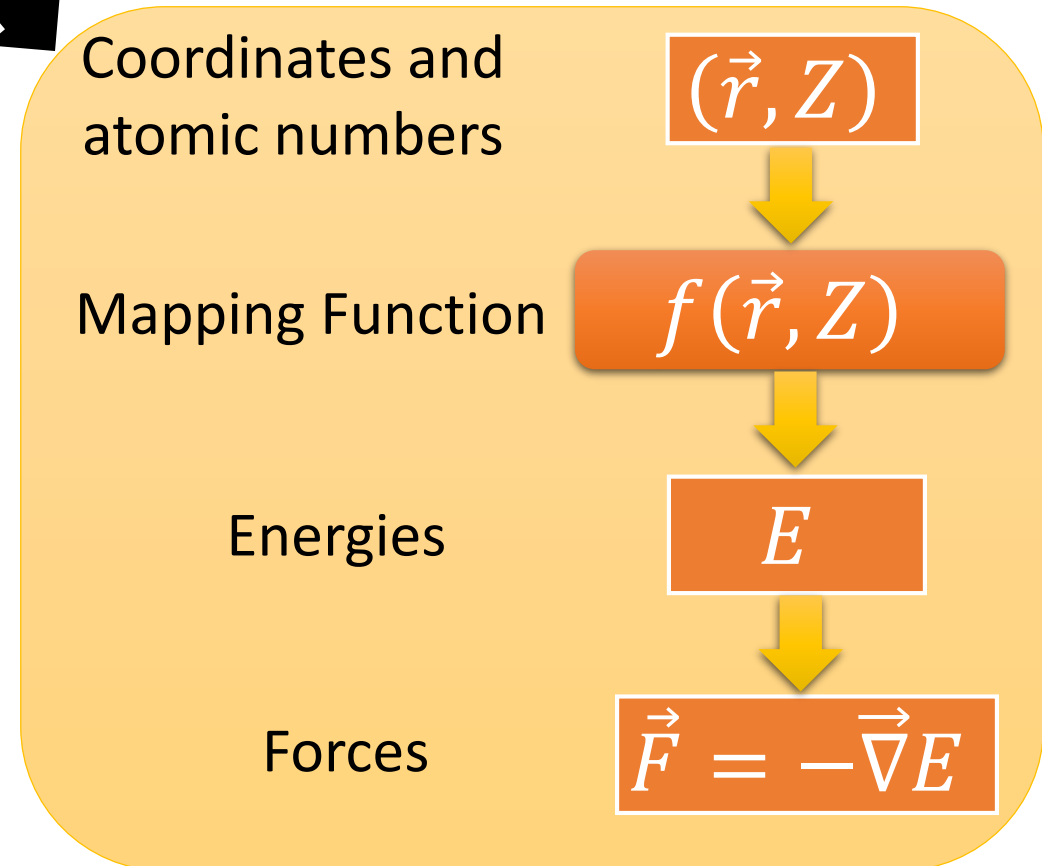




# Design principles for ML potentials

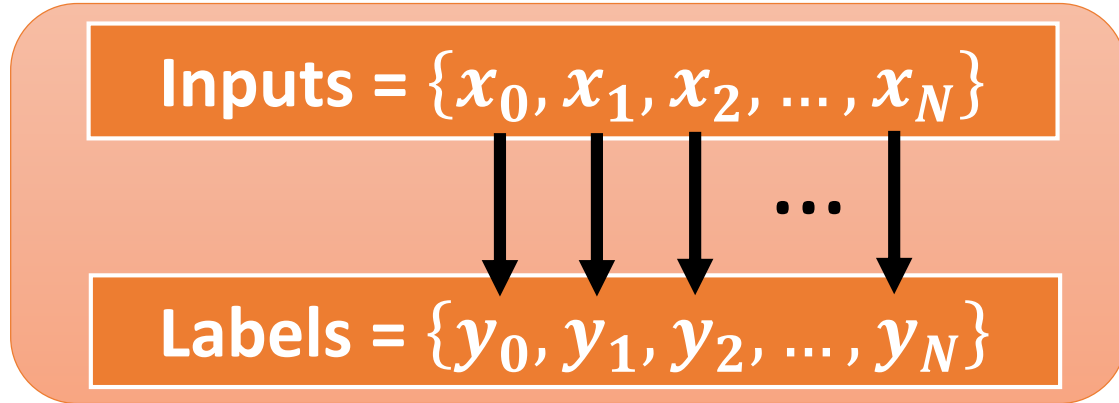
Mapping from coordinates  $R \rightarrow$  Energy (& Forces) but with no a-priori functional form

- Fast, accurate, and reproducible
- Reactive
- No “atom typing” required
- Conserves energy (in MD)
- Extensible to new, larger systems of atoms
- Highly automated parametrization
- Systematically improvable



# Machine learning basics

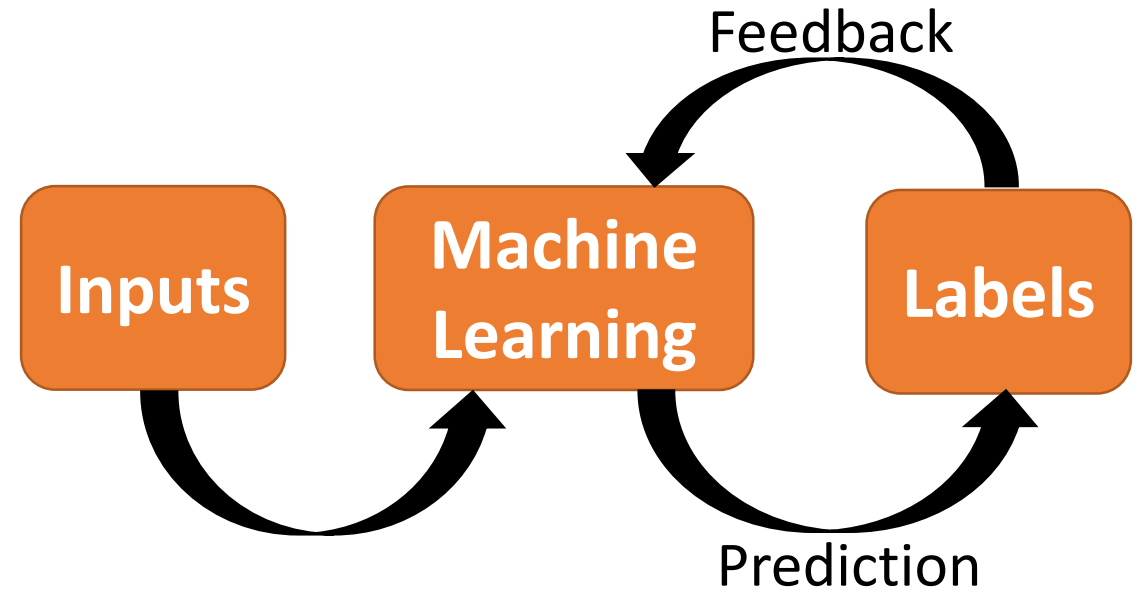
## Training dataset for supervised learning



## Types of tasks

- Regression
- Classification

## Supervised Learning



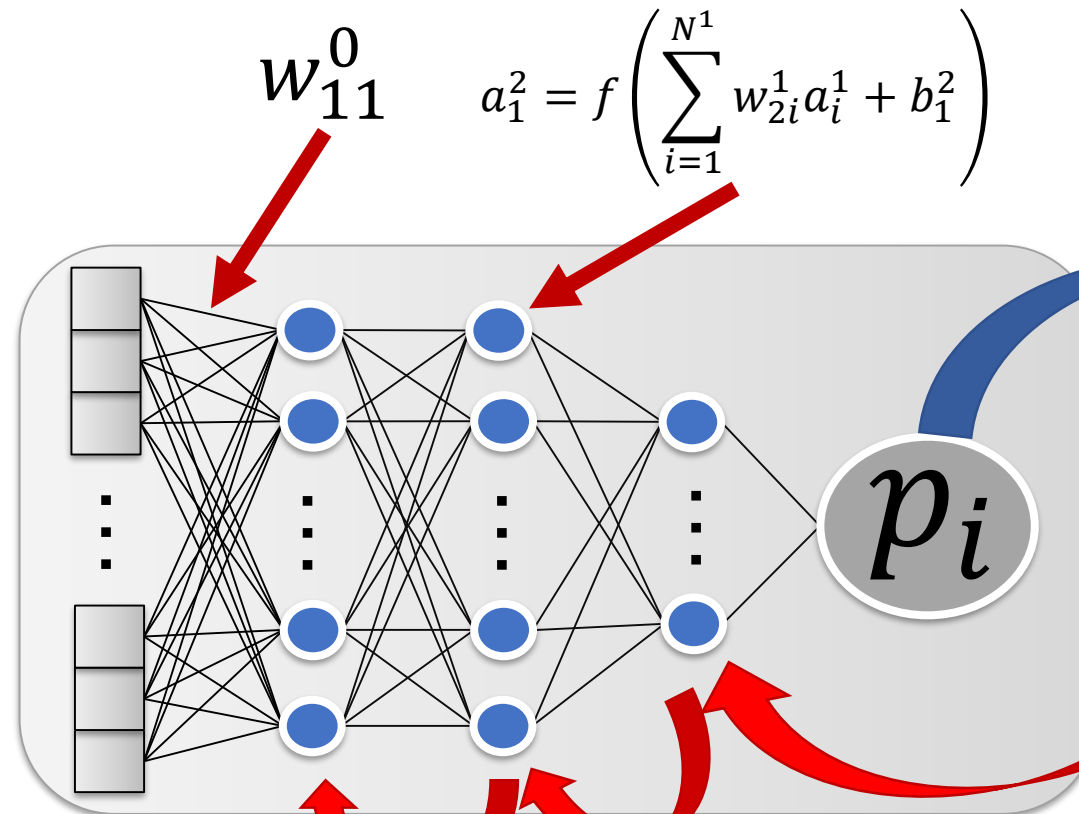
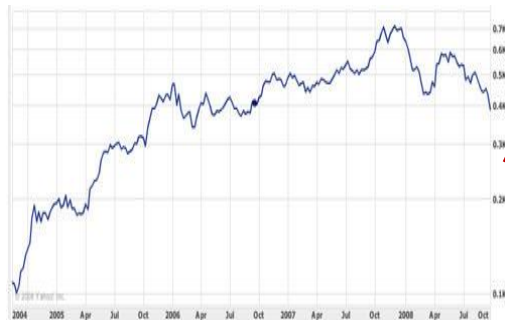
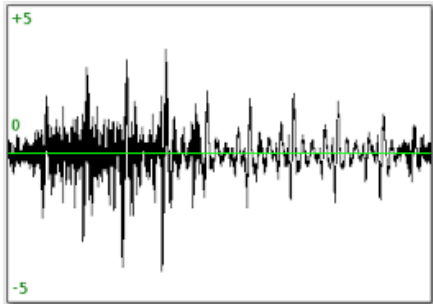
## A few applications

- Image recognition
- Social media moderation
- Stock market prediction

# Deep Learning

neural networks (NN)

1 1 5 4 3  
7 5 3 5 3  
5 5 9 0 6  
3 5 2 0 0

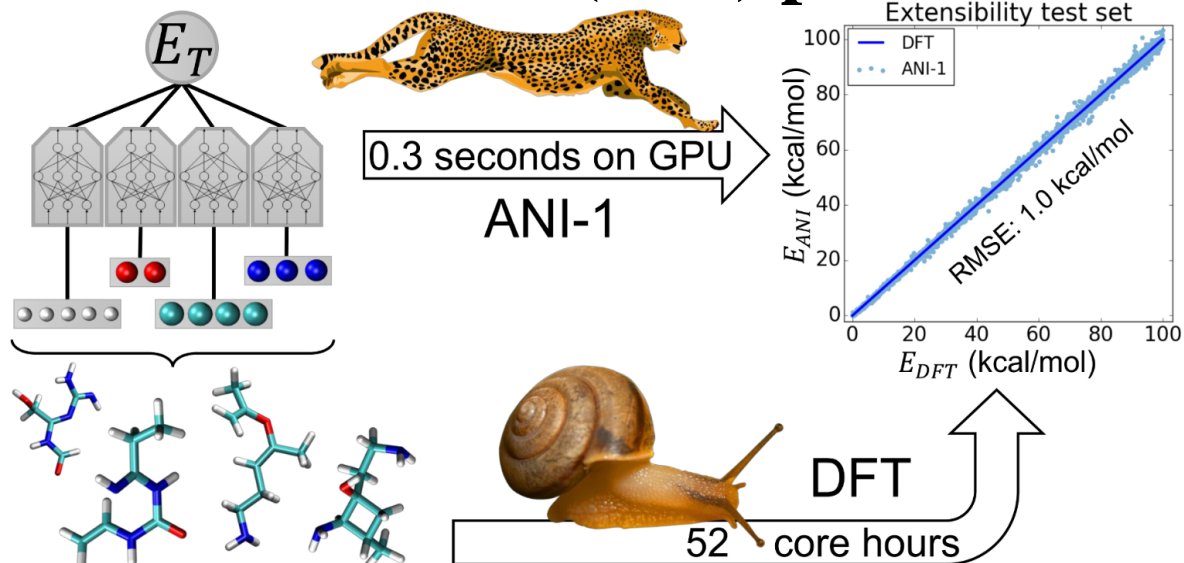


Back Propagation to compute  $\frac{\partial C}{\partial w_{ki}^j}$   
and update weights with SGD

$$f(x) = \tanh(x)$$

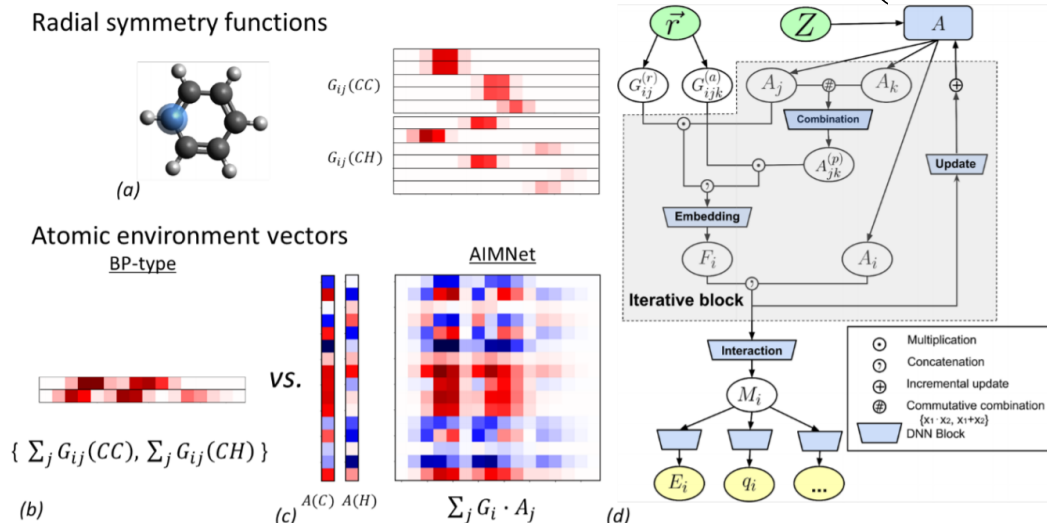
$$C = \frac{1}{M} \sum_j^M (p_i - \widehat{A}_i)^2$$

# ANAKIN-ME (ANI) potentials



JS Smith, O Isayev, AE Roitberg, *Chem. Sci.*, 2017, **8**, 3192-3203

## Atoms-in-molecule neural network (AIM-Net)

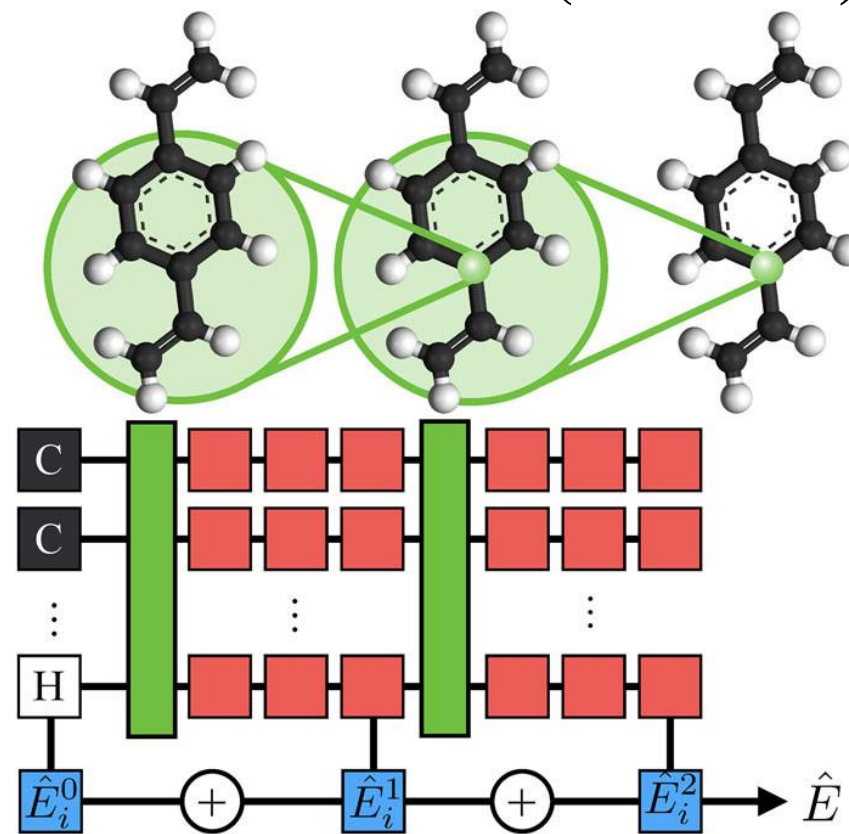


R Zubatyuk, JS Smith, J Leszczynski, O Isayev

<https://doi.org/10.26434/chemrxiv.7151435.v2> 2018

# Our work on developing ML potentials

## Hierarchal Interacting Particle Neural Networks (HIP-NN)



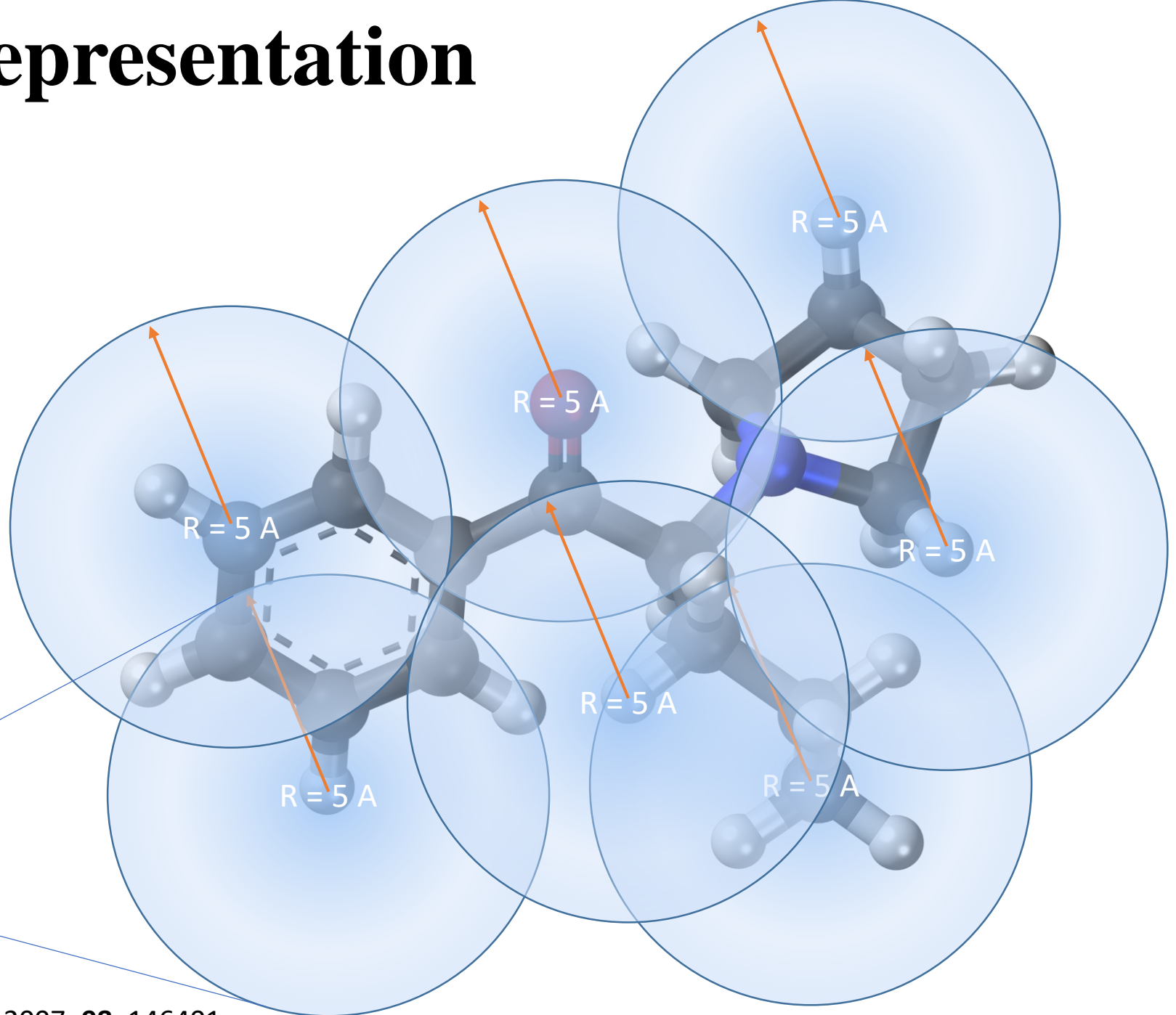
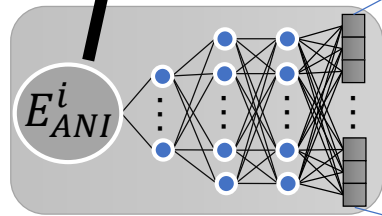
N Lubbers, JS Smith, K Barros

*J. Chem. Phys.*, 2018, **148**, 241715

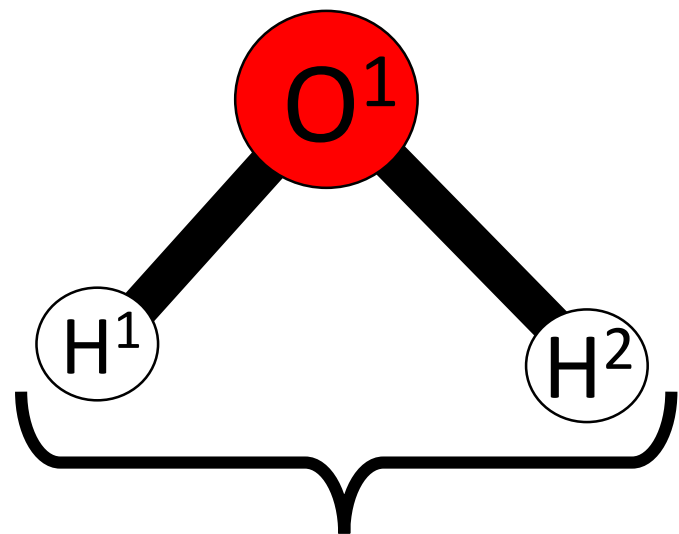
# Molecular Representation

Energy is given by  
a sum over atomic  
contributions

$$E = \sum_i^N E_i$$

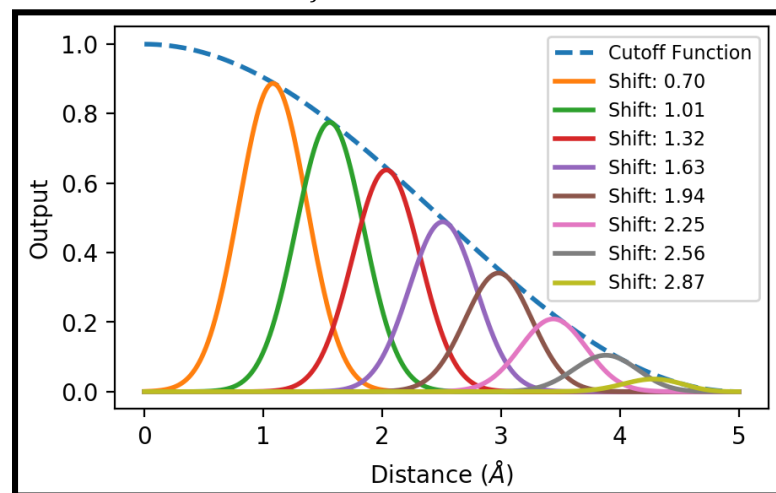


# Descriptors for the ANI ML-based potential



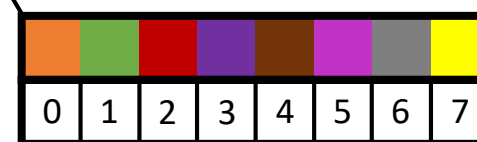
## Radial Descriptors

$$G_m^R = \sum_{j \neq i}^{\text{All Atoms}} e^{-\eta(R_{ij}-R_s)^2} f_c(R_{ij})$$



## Cutoff function

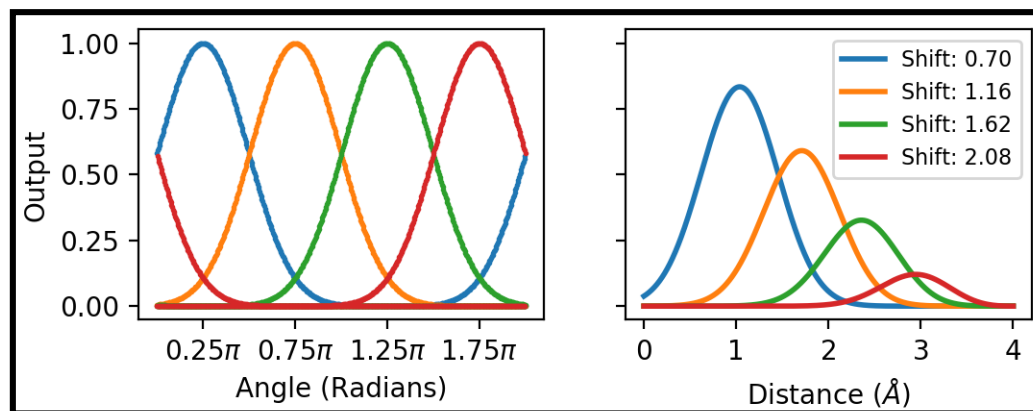
$$f_c(R_{ij}) = \begin{cases} 0.5 \times \cos\left(\frac{\pi R_{ij}}{R_c}\right) + 0.5 & \text{for } R_{ij} \leq R_c \\ 0.0 & \text{for } R_{ij} > R_c \end{cases}$$



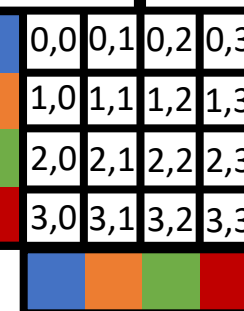
Concatenate

## Angular descriptors

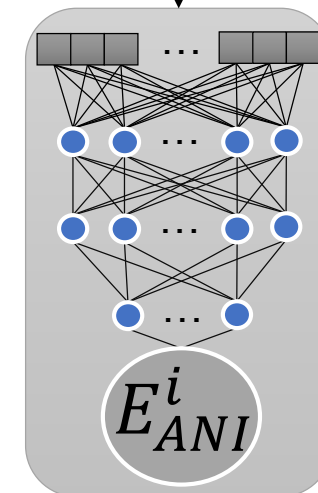
$$G_m^{A_{\text{mod}}} = 2^{1-\zeta} \sum_{j,k \neq i}^{\text{All Atoms}} (1 + \cos(\theta_{ijk} - \theta_s))^\zeta \exp\left[-\eta\left(\frac{R_{ij}^2 + R_{ik}^2}{2} - R_s\right)^2\right] f_c(R_{ij}) f_c(R_{ik})$$



Angular



Radial



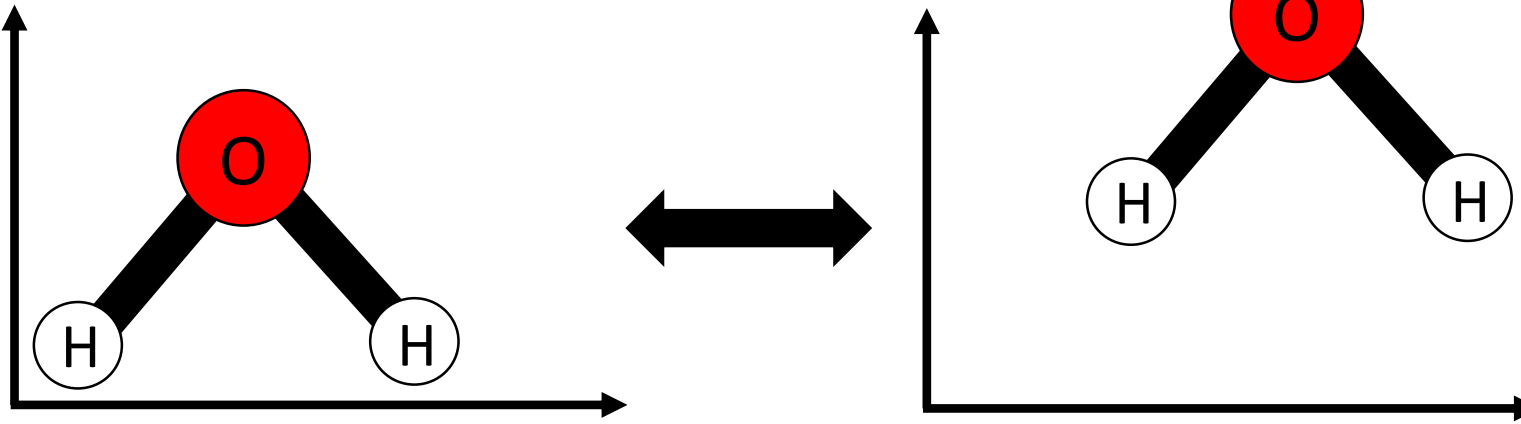
O <sup>1</sup>	H <sup>1</sup>	H <sup>2</sup>
D(H <sup>1</sup> )	D(H <sup>2</sup> )	D(H <sup>1</sup> )
D(H <sup>2</sup> )	D(O <sup>1</sup> )	D(O <sup>1</sup> )
A(H <sup>1</sup> ;H <sup>2</sup> )	A(O <sup>1</sup> ;H <sup>2</sup> )	A(O <sup>1</sup> ;H <sup>1</sup> )

D = distance to atom

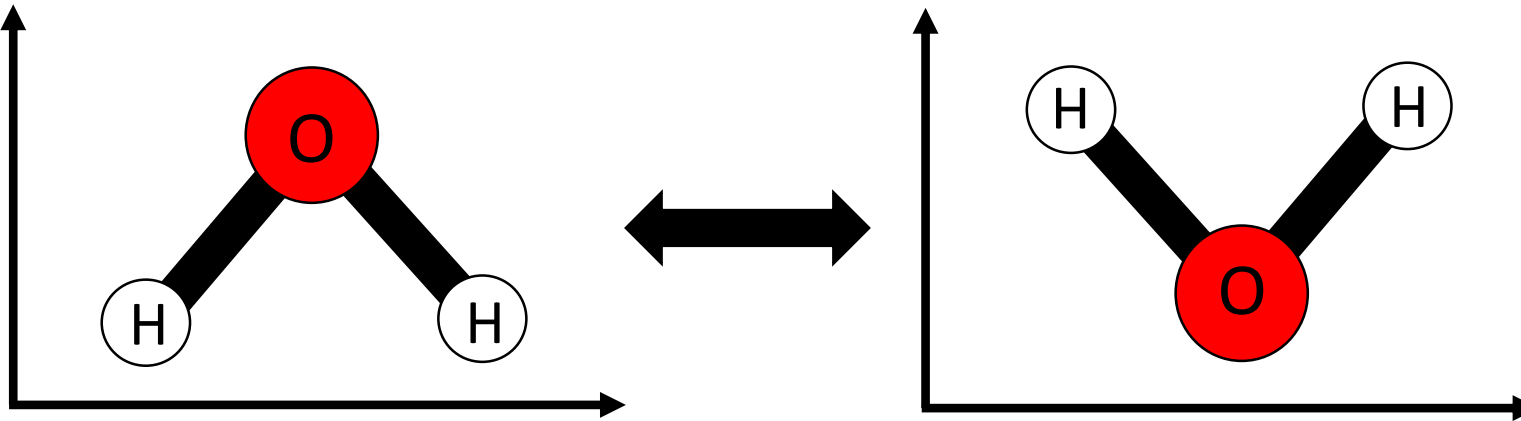
A = Angle between atoms

**Energy must be invariant  
with respect to:**

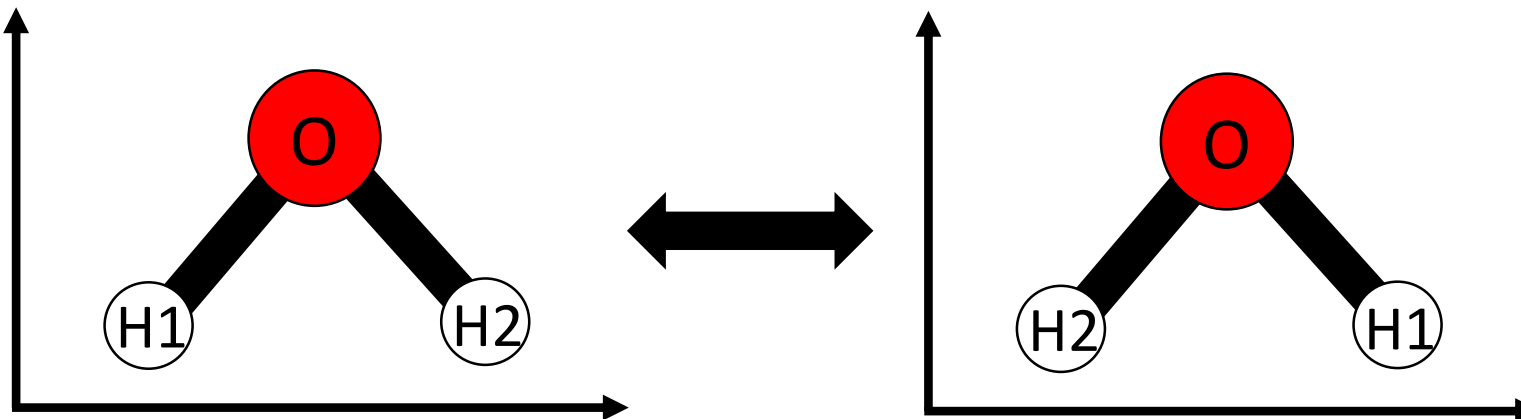
Translation



Rotation

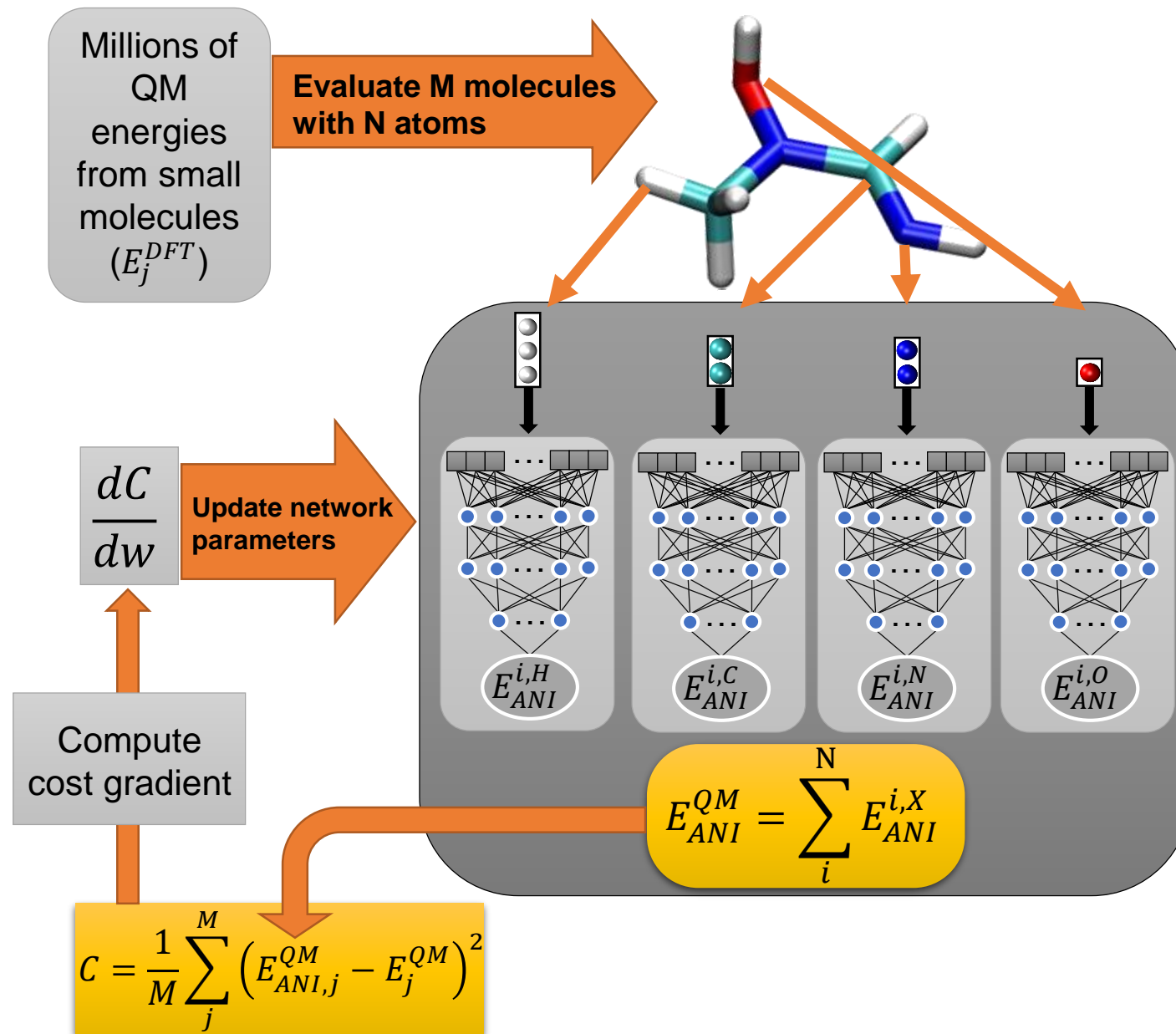


Permutation

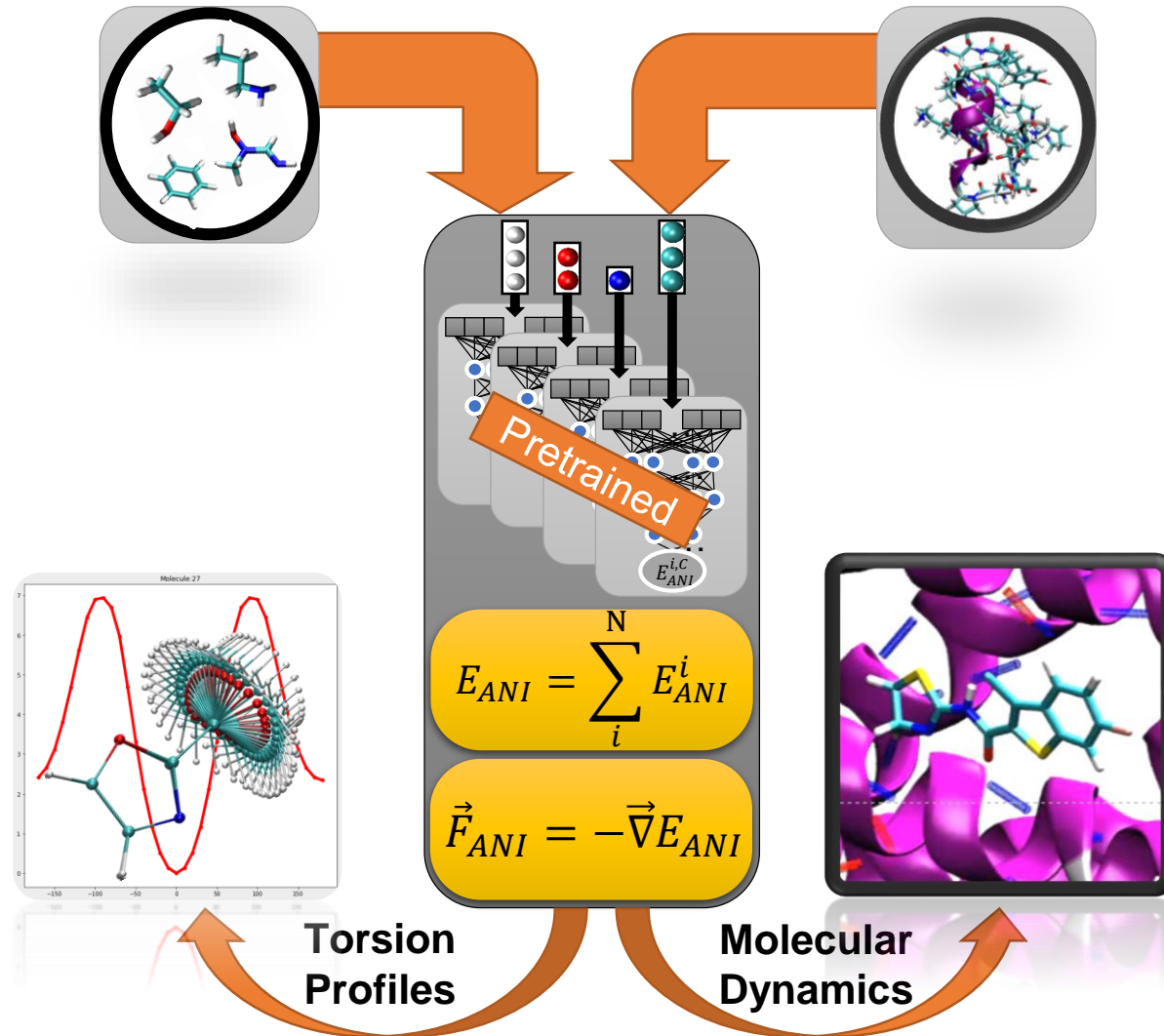




# ANI potential training



# ANI potential application

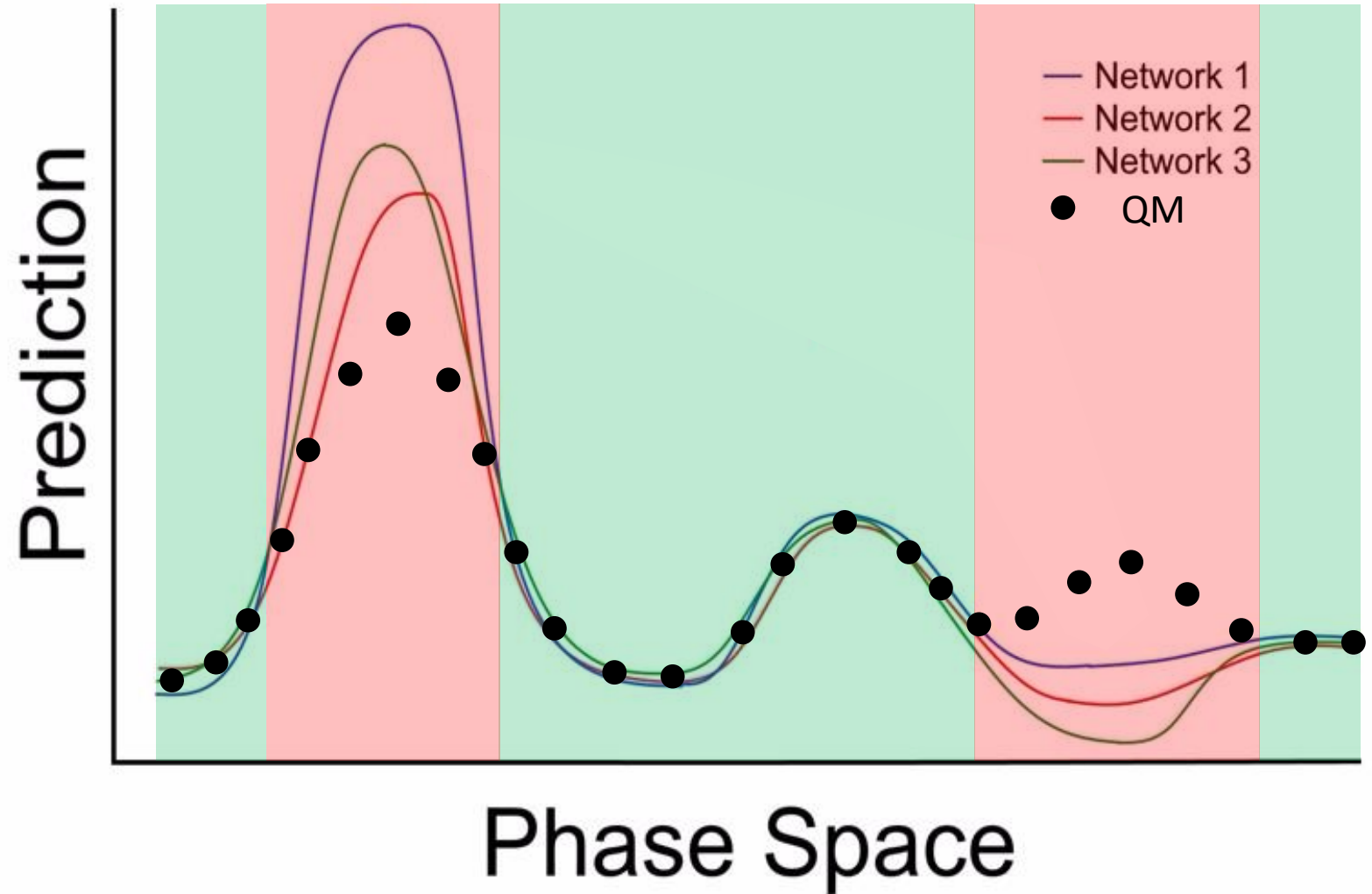


# Can we predict when the model is wrong?

Ensemble  
disagreement  
can drive  
data  
generation

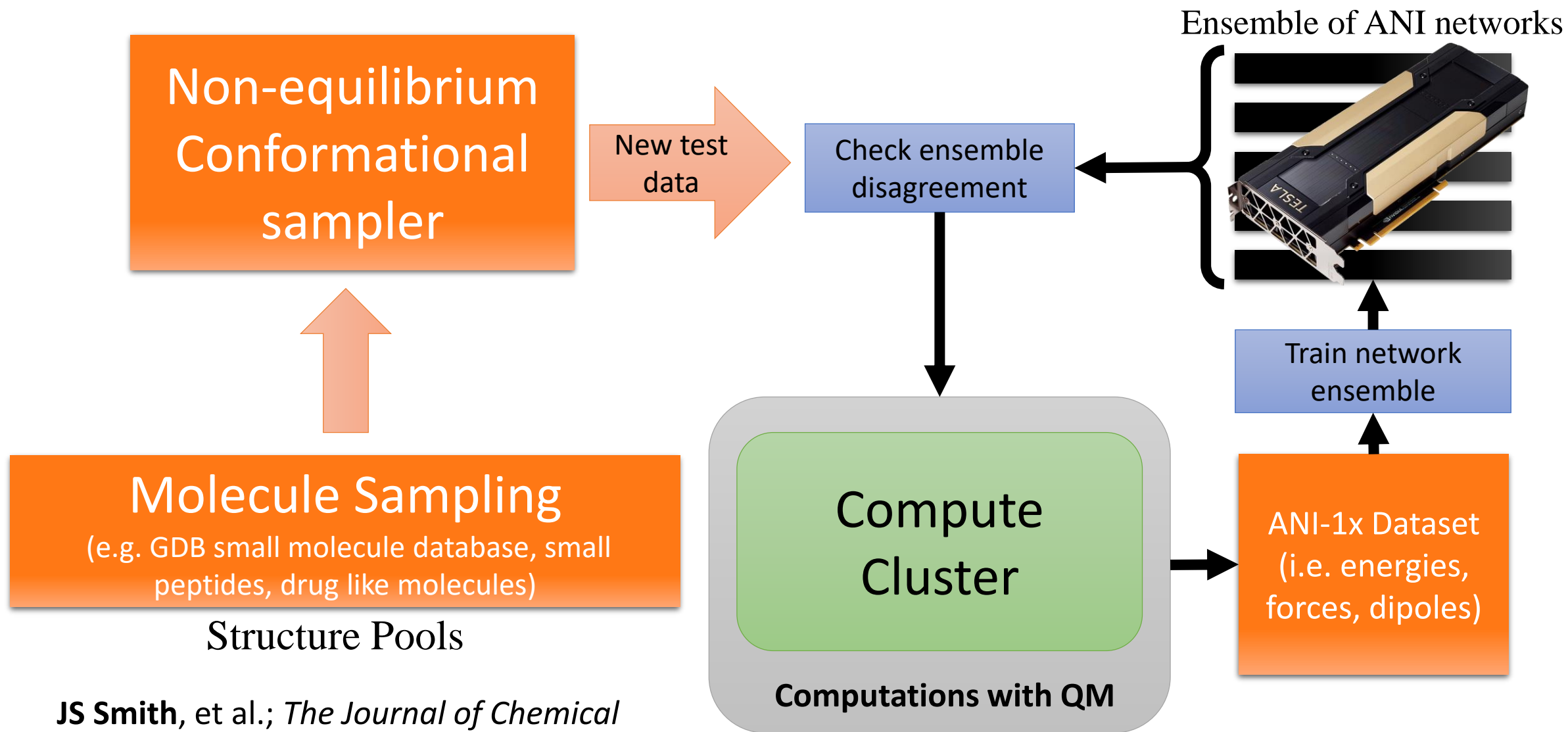
Good data  
coverage

Bad data  
coverage



# Active Learning - The Big Picture

An automated and self-consistent data generation framework

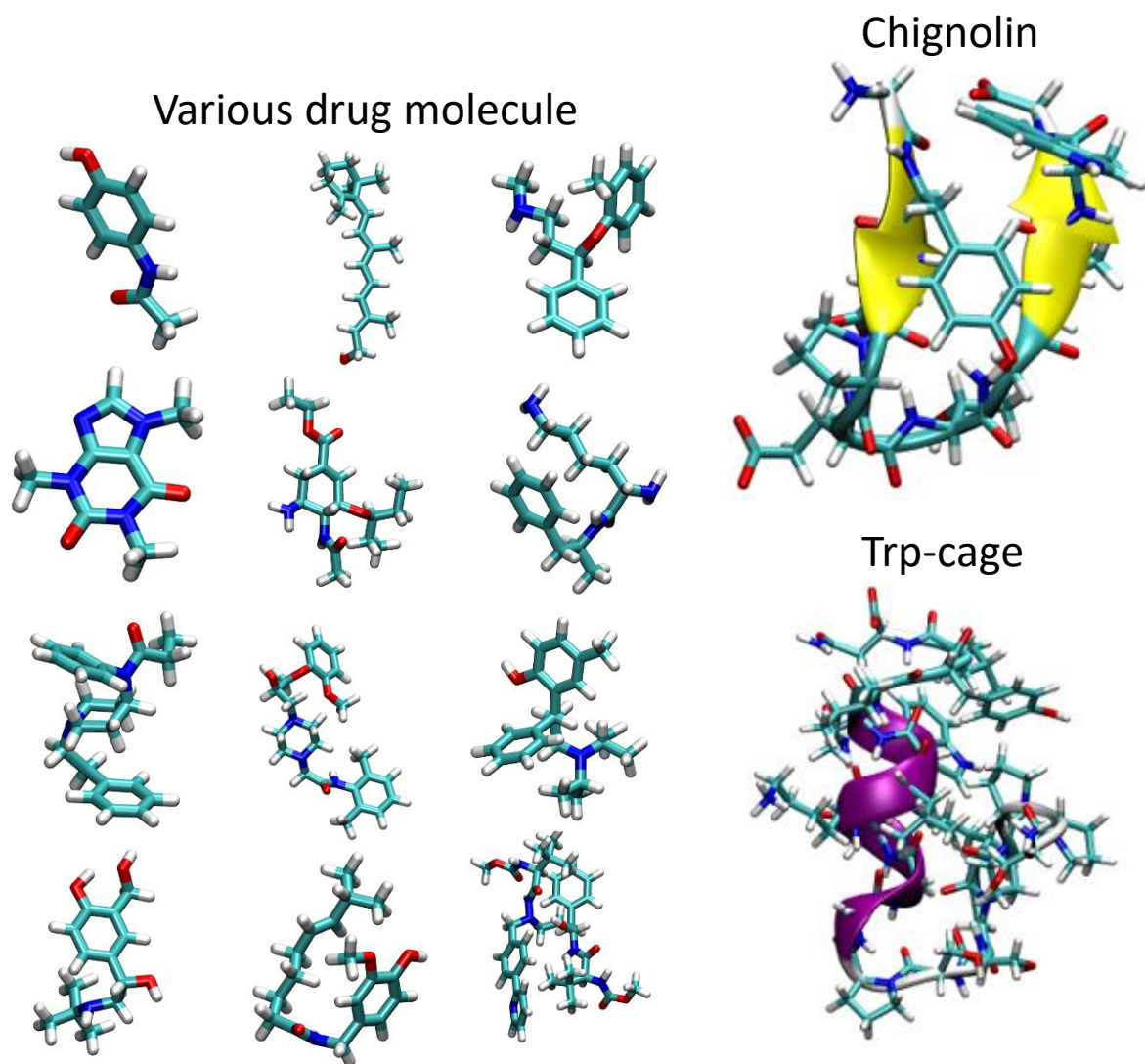


JS Smith, et al.; *The Journal of Chemical Physics*, (2018), 148 (24), 241733

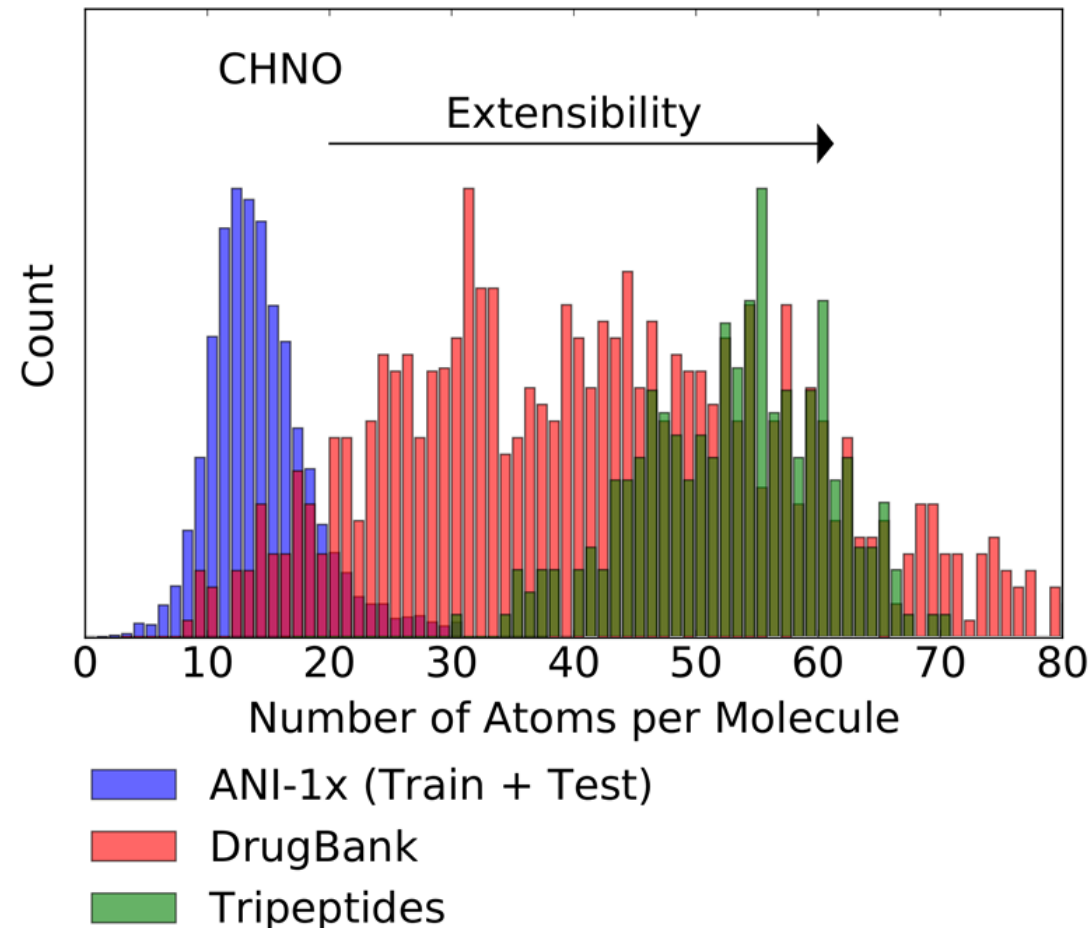
# Testing transferability and extensibility

## ANI-MD Benchmark

128 frames from 1ns trajectories @ 300K for each:



## DrugBank and Tripeptide Benchmarks



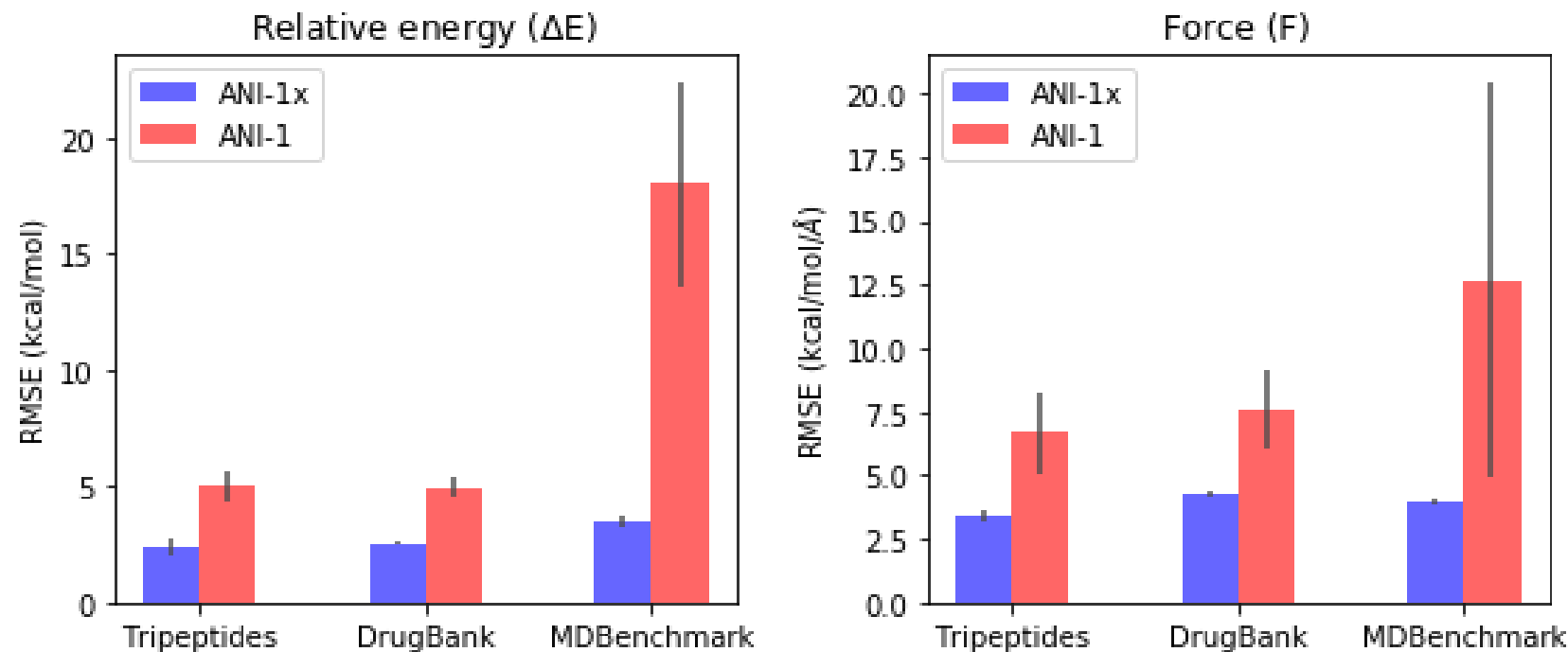
**JS Smith**, et al.; *The Journal of Chemical Physics*, (2018), 148 (24), 241733

# Active-learning results vs. random sampling

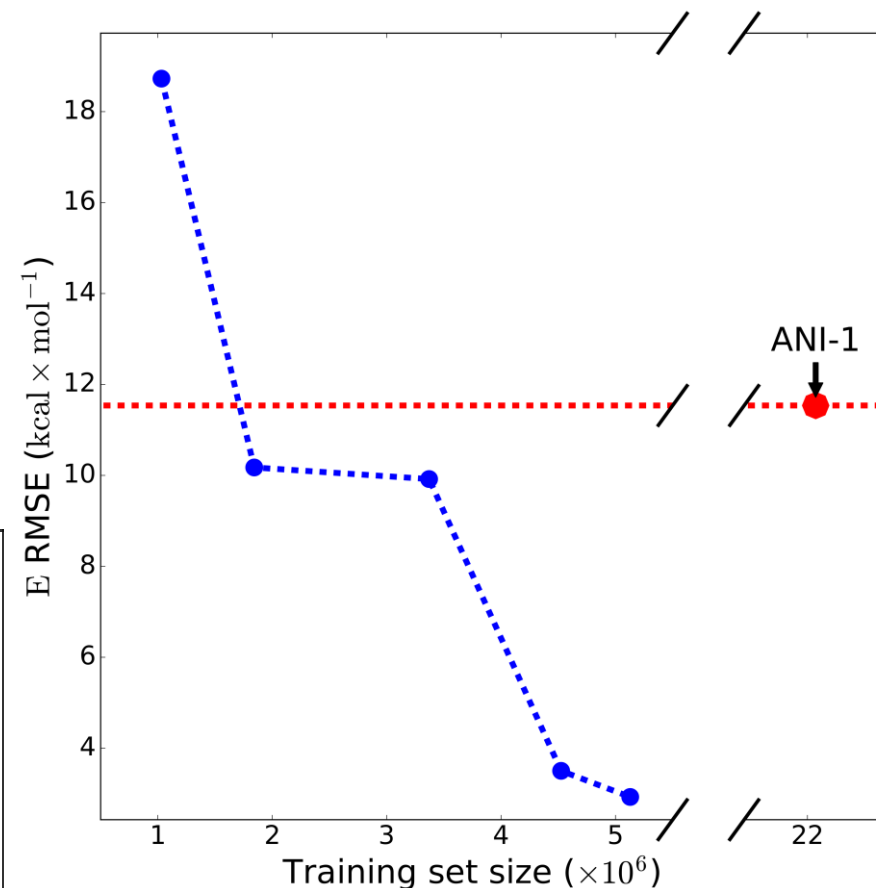
## Dataset size comparison

	ANI-1	ANI-1x
Datapoints	22M	5M

## Relative E and F RMSE comparison



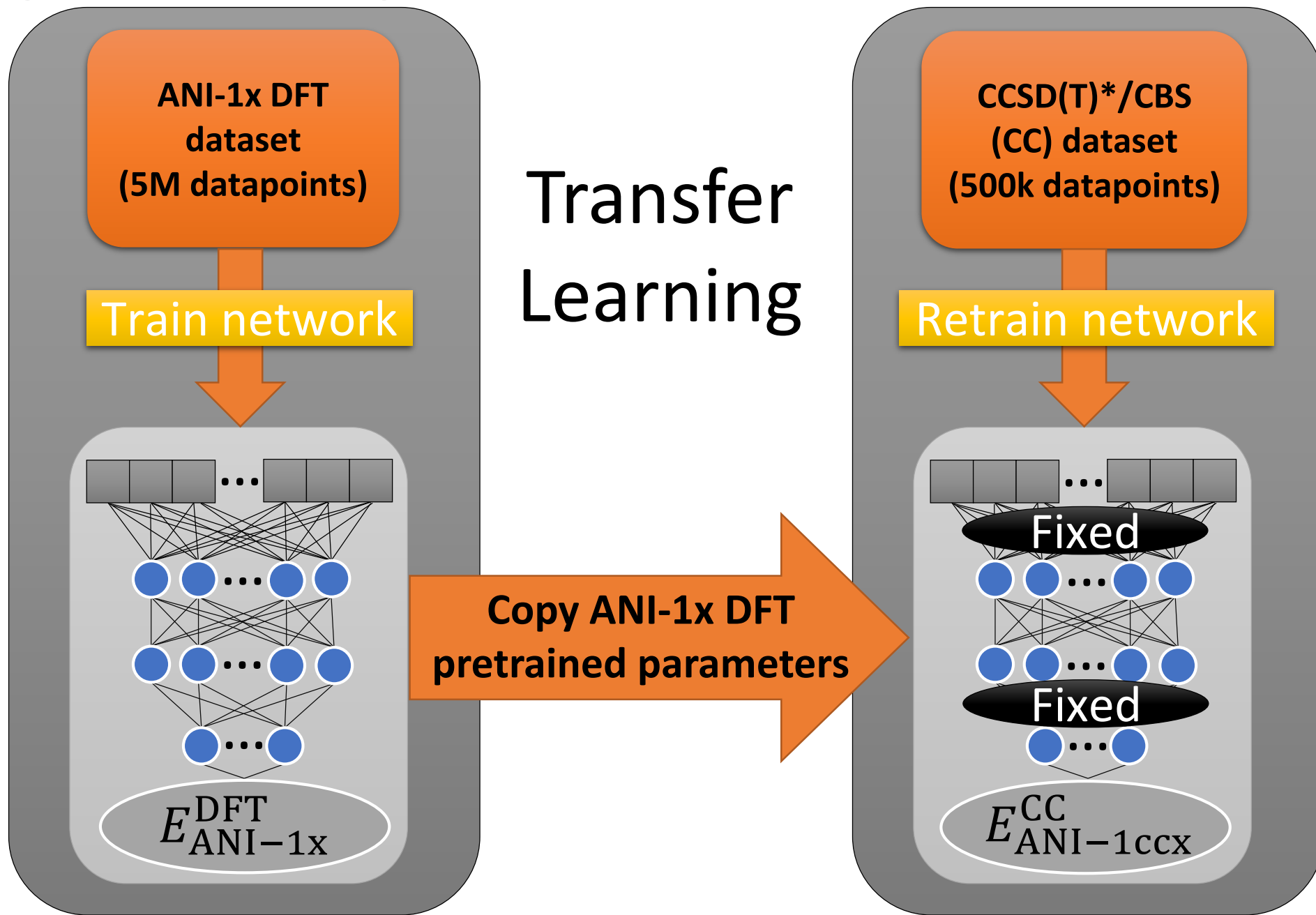
## Active learning progression



**JS Smith, et al.; *The Journal of Chemical Physics*, (2018), 148 (24), 241733**

# Transferring knowledge from DFT to CCSD(T)

- Subsample 10% of ANI-1x training data (0.5M of 5M)
- Recompute CCSD(T)/CBS level
- 340k parameters fixed, re-train 60k
- $10^7$  faster than DFT

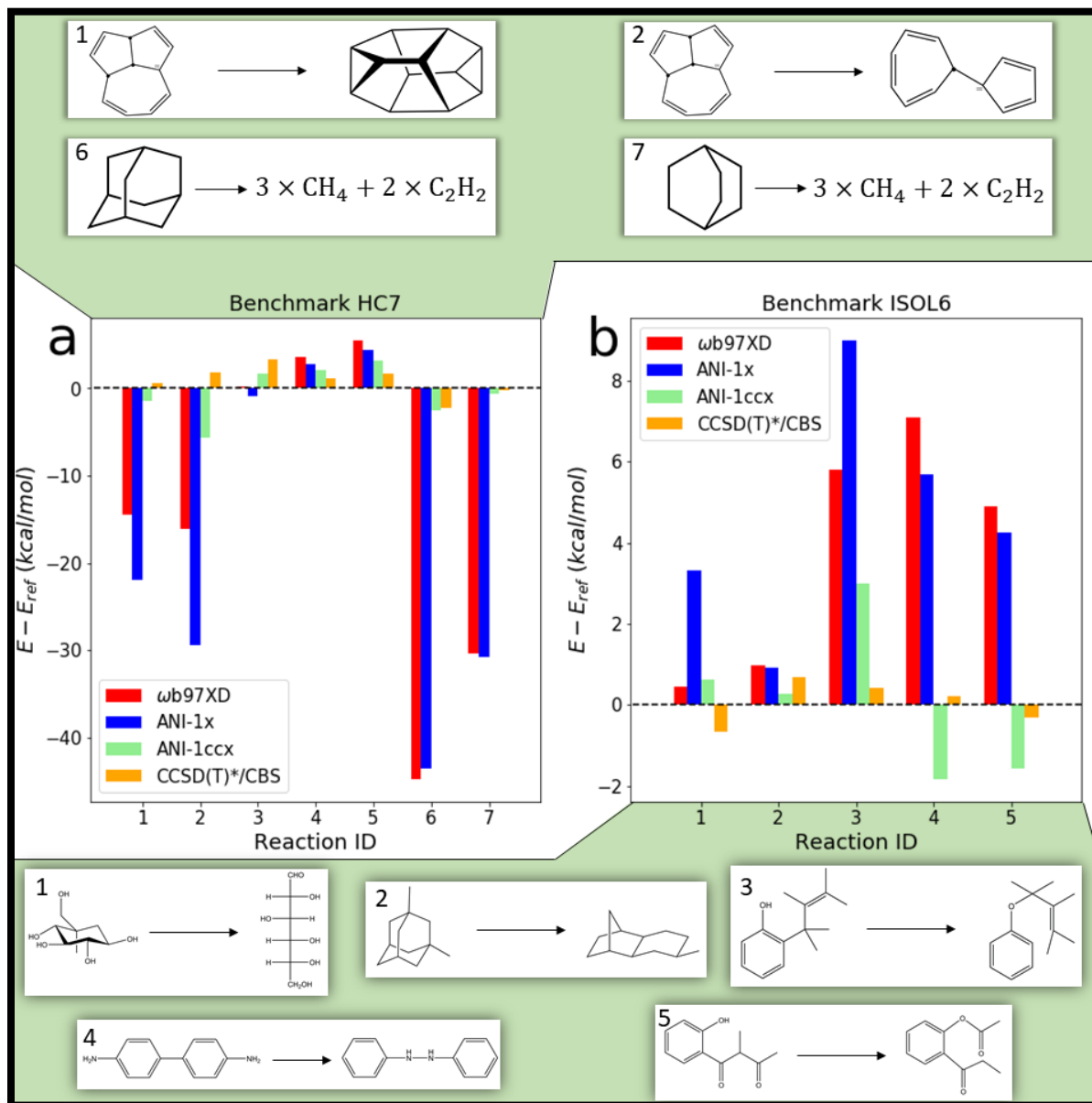




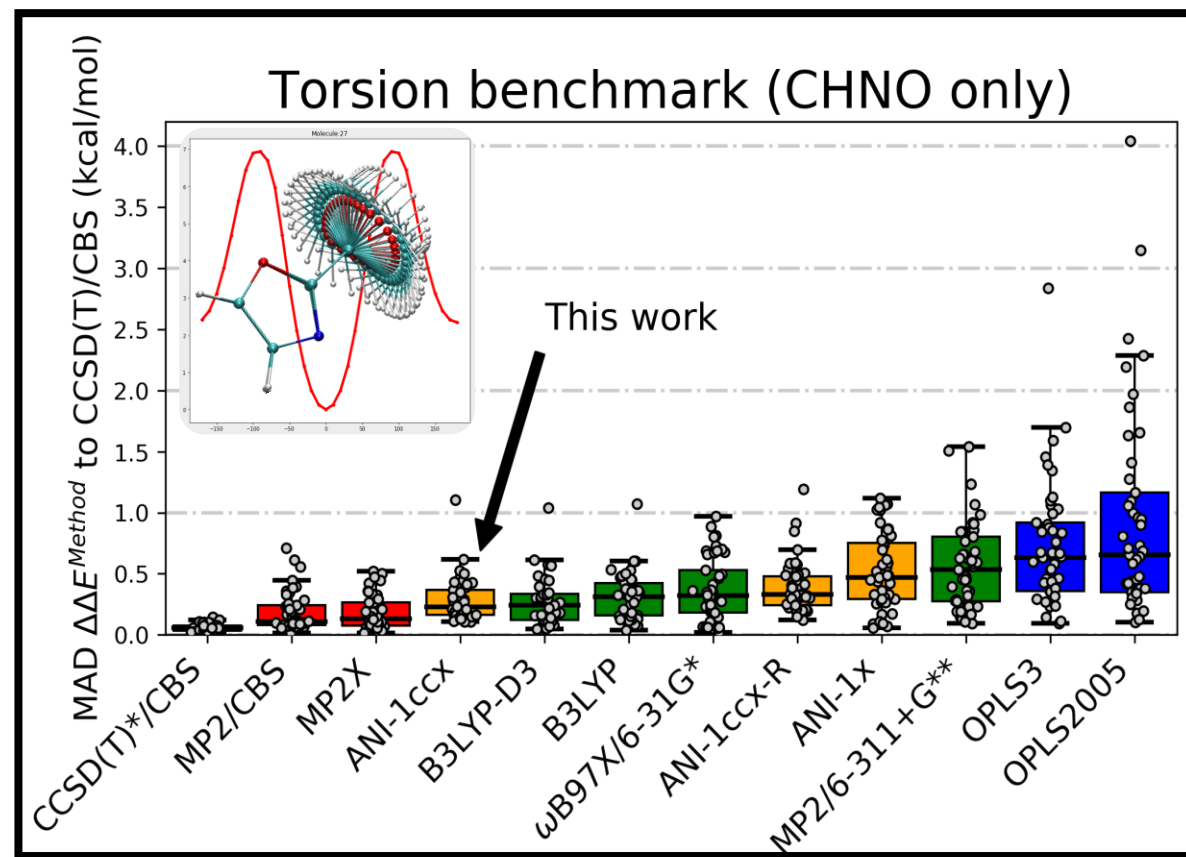
# Outsmarting Quantum Chemistry Through Transfer Learning

JS Smith, B Nebgen, R Zubatyuk, N Lubbers, C Devereux, K Barros, S Tretiak, O Isayev, A Roitberg

<https://doi.org/10.26434/chemrxiv.6744440.v1> **2018** (under review at Nat. Comm.)



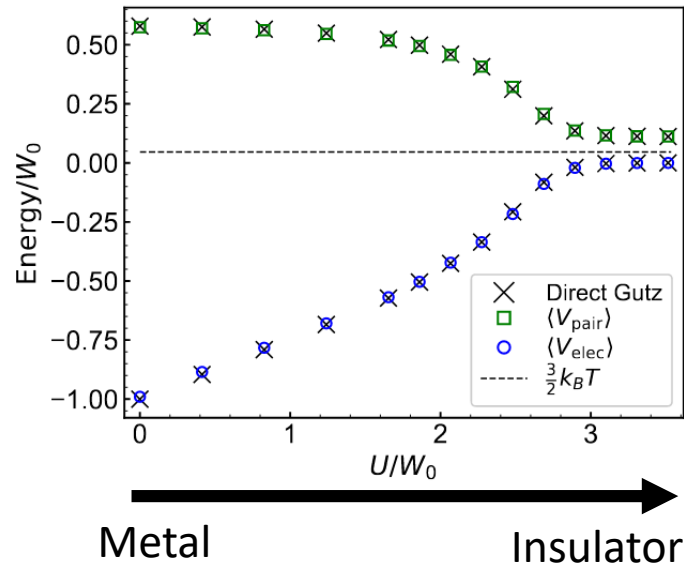
- New ANI-1ccx model outperforms DFT on reaction energies and torsional profiles
- A 24 core hours calculation for CCSD(T)/CBS takes 2 GPU microseconds for ANI-1ccx



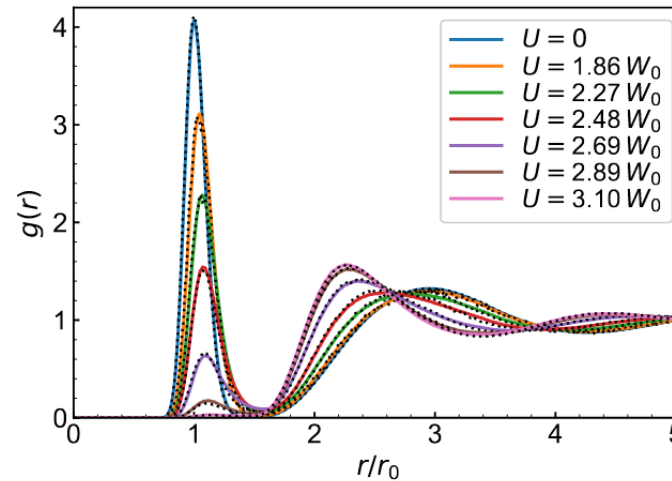
# Machine learning for molecular dynamics with strongly correlated electrons

Hidemaro Suwa, Justin S. Smith, Nicholas Lubbers, Cristian D. Batista, Gia-Wei Chern, Kipton Barros

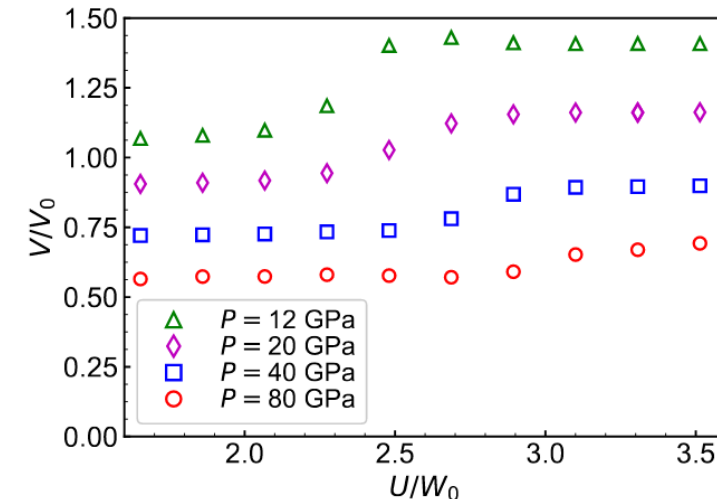
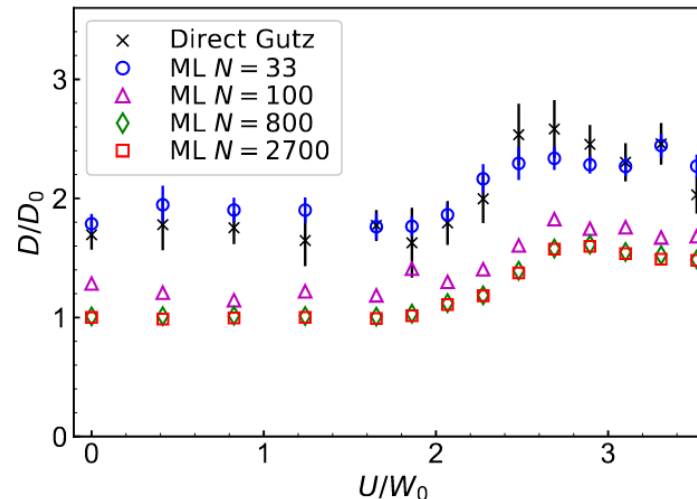
<https://arxiv.org/abs/1811.01914> **2018** (under review at Phys. Rev. Lett.)



Accurately reproduces a Mott transition on systems of 2700 atoms.



Trained an ML model to a toy system with variable Hubbard  $U$  through the Gutzwiller approximation



# ML metal potentials with active learning!

Our approach: minimize use of expert knowledge for **maximum generality**

## Los Alamos Team

Benjamin Nebgen – T-1

Kipton Barros – T-1

Saryu Fensin – MST-8

Tim Germann – T-1

Leonid Burakovsky – T-1

Nicholas Lubbers – CCS-3

Sergei Tretiak – T-1

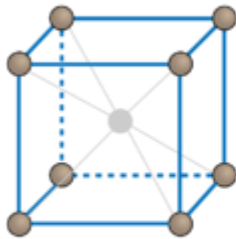
$$\begin{array}{ll} \text{Energy} & E[\mathbf{r}_1, \mathbf{r}_2, \dots] \\ \text{Force} & \mathbf{f}_i = -\nabla_i E \end{array}$$

Applications

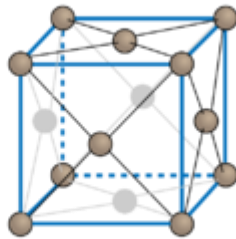
$$m \frac{d^2 \mathbf{r}_i}{dt^2} = \mathbf{f}_i$$

ML Potential

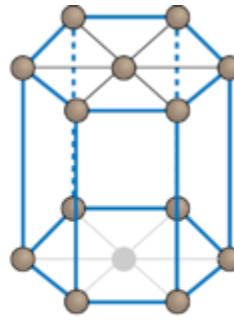
Crystal structures



**Cubic body centered (bcc)**  
Fe, V, Nb, Cr



**Cubic face centered (fcc)**  
Al, Ni, Ag, Cu, Au



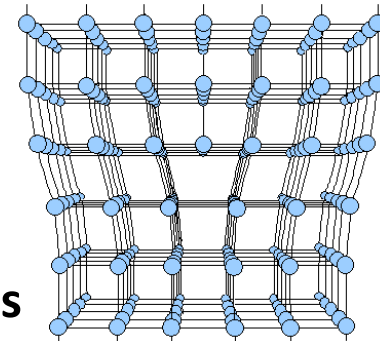
**Hexagonal**  
Ti, Zn, Mg, Cd

Melts

Amorphous solid

Extreme conditions

Defects

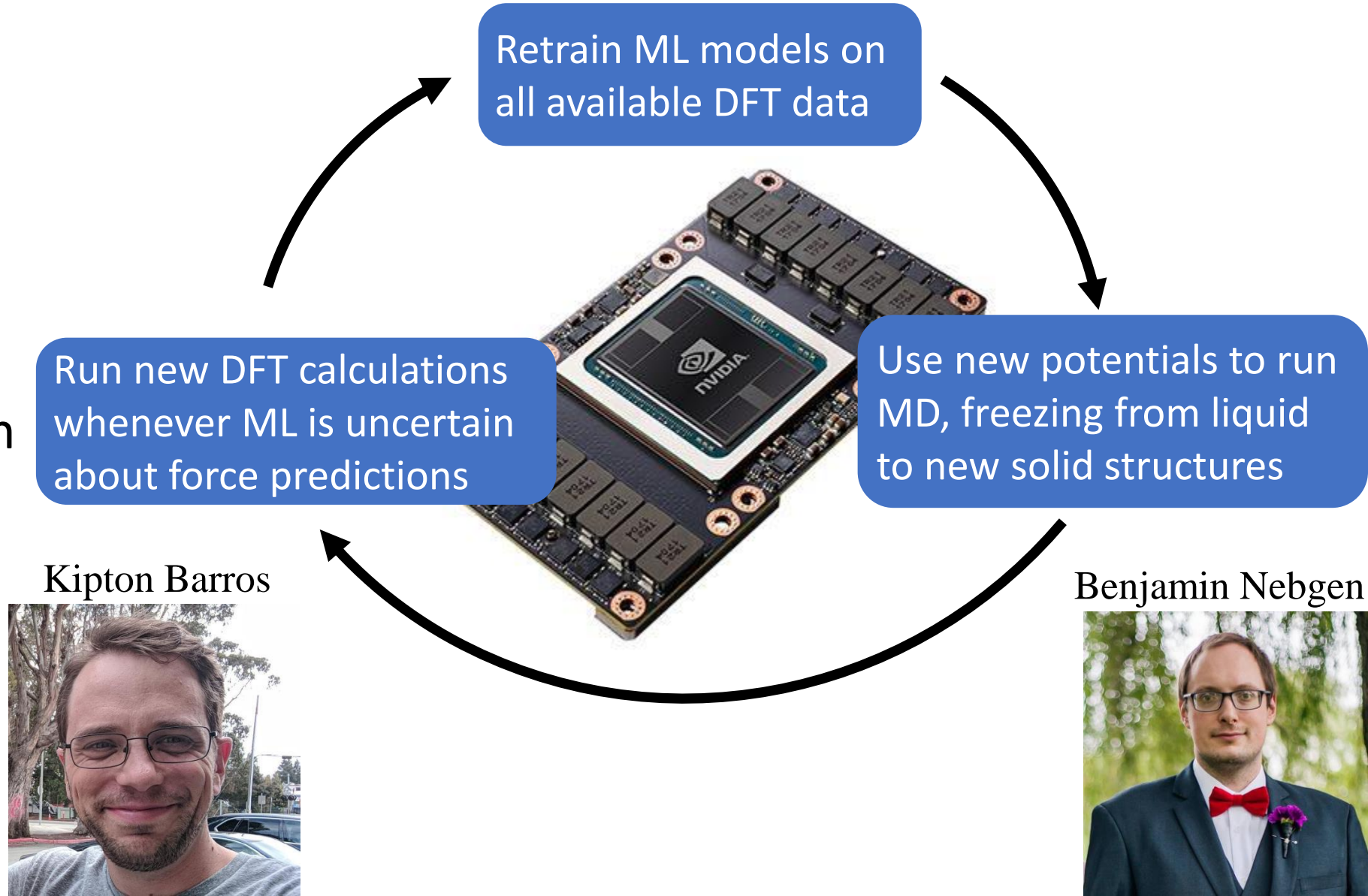


# The General AL Framework for ML potentials

A codebase for active learning on large GPU clusters (e.g. Sierra or Summit)

## Design Principles:

- Fully autonomous
- Interchangeable QM
- Interchangeable ML
- Assortment of built-in sampling methods
- Built-in testing suite
- Capable of scaling to 1000s of nodes

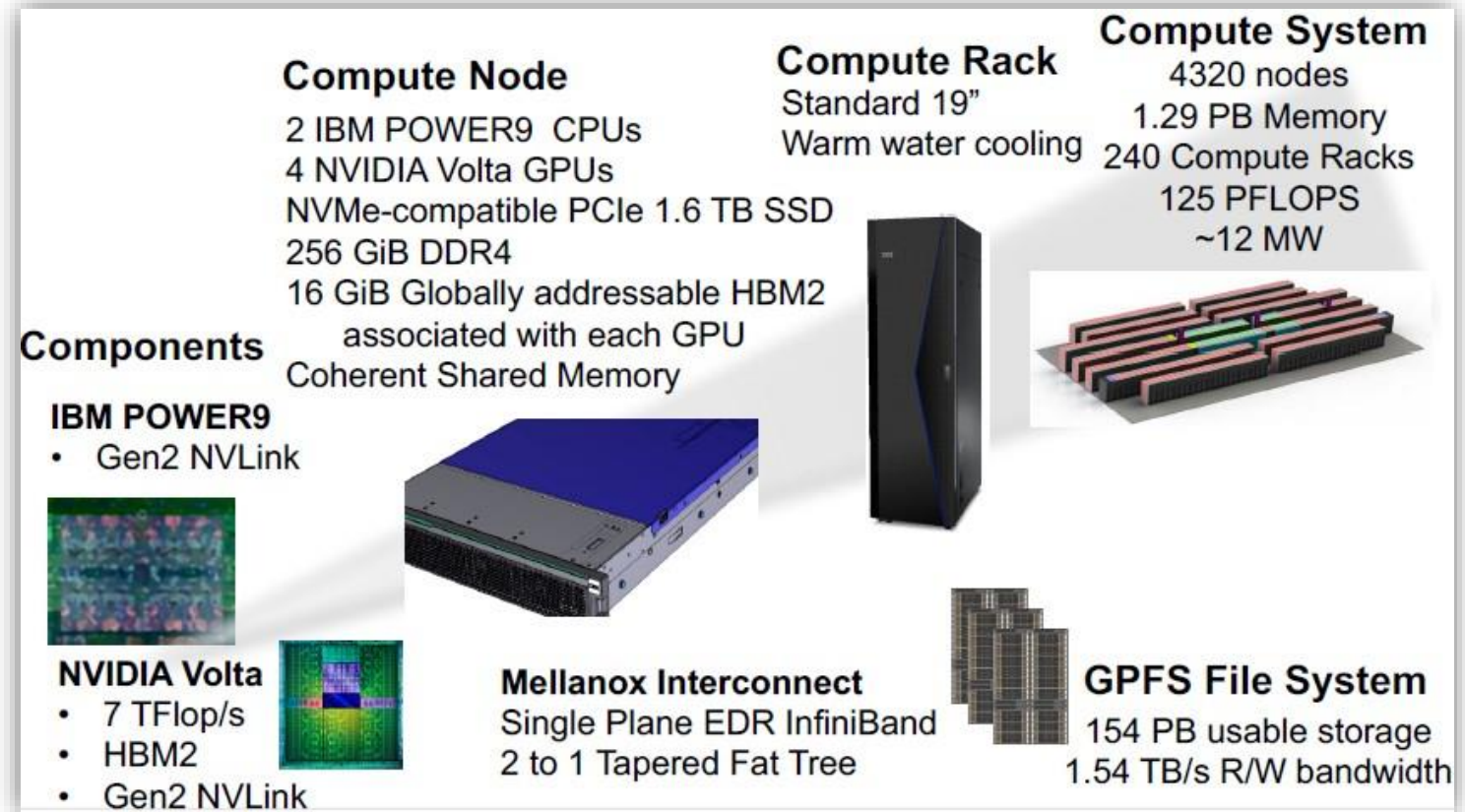




## Application of active learning to build ML potentials

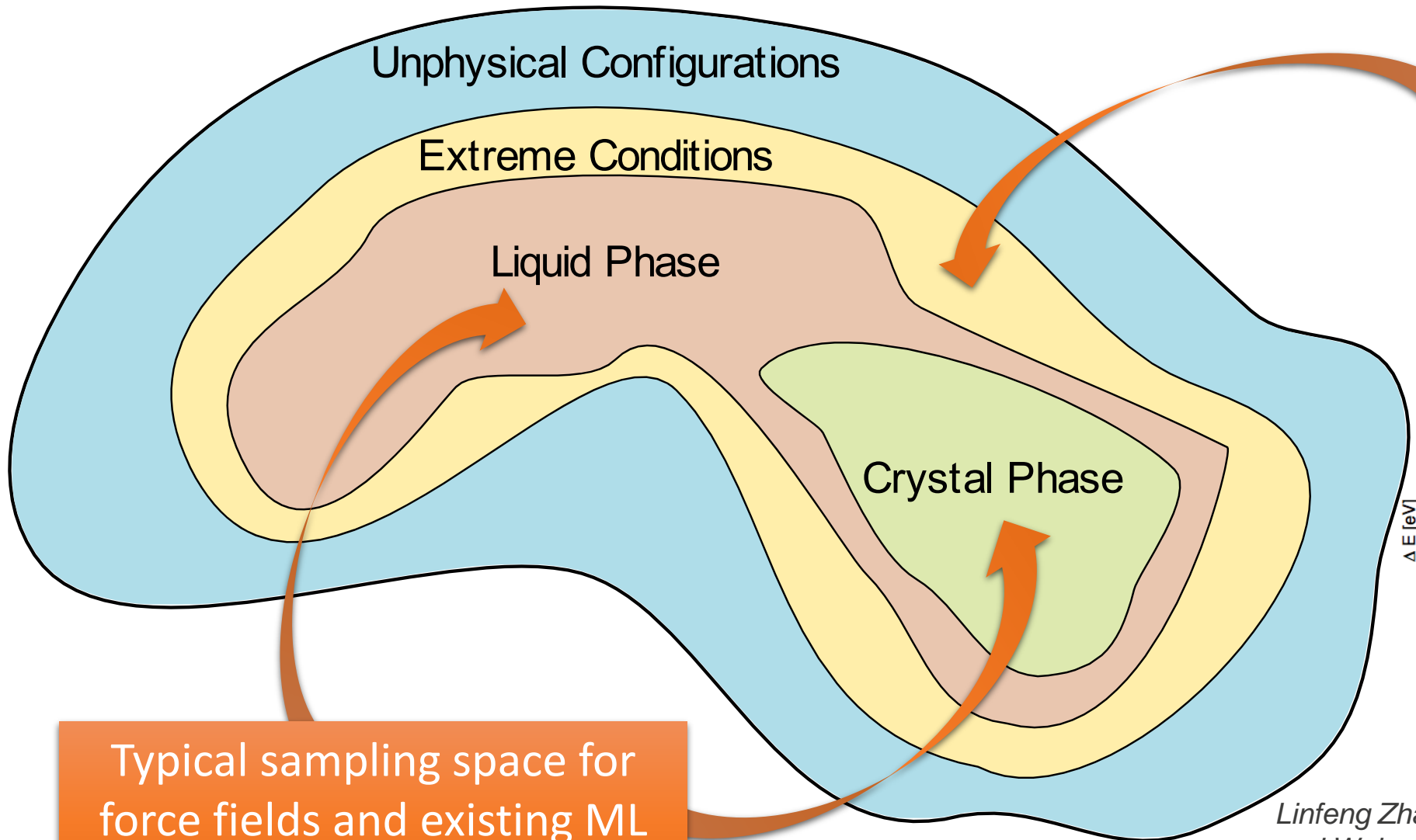
- Build a general active learning framework
- Framework interface with Quantum Espresso (QE) for DFT
- In 2 months we ran 10-20k DFT calculations on systems with 50-200 metal atoms
- Elements explored **Al**, **Sn**, Ga, Cu
- We are still evaluating results

## Open science access on LLNL's Sierra super computer



# How should we sample to build a general model?

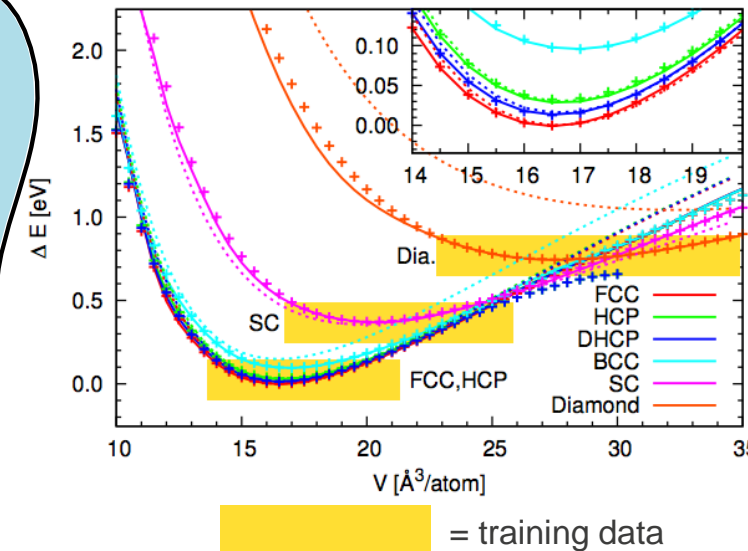
All possible configurations for a metal



Where extreme conditions and rare events exist (e.g. shock simulations)

Typical sampling space for force fields and existing ML potentials

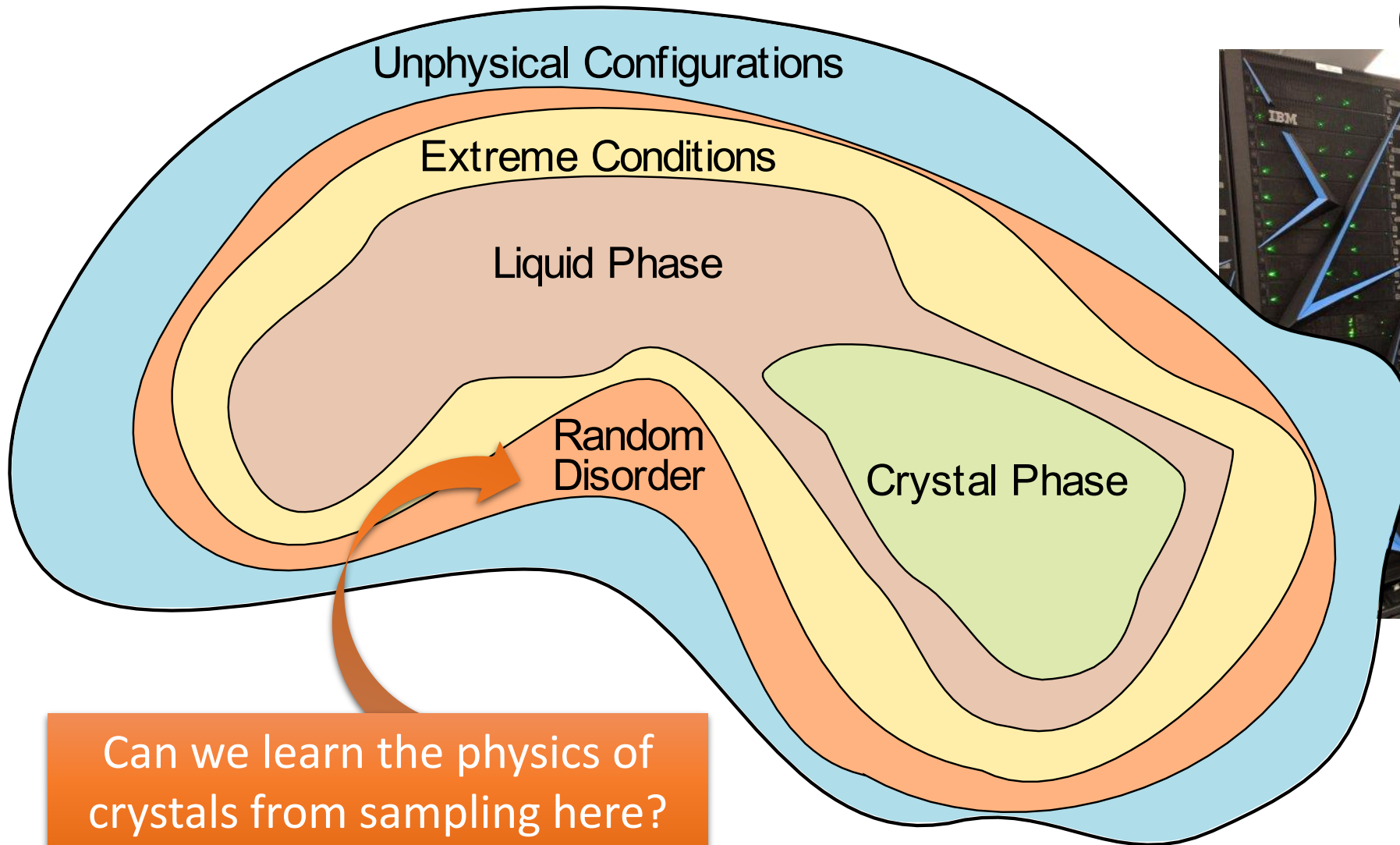
Recent published work



Linfeng Zhang, De-Ye Lin, Han Wang, Roberto Car, and Weinan E, Active Learning of Uniformly Accurate Inter-atomic Potentials for Materials Simulation, [arXiv:1810.11890]

# How should we sample to build a general model?

All possible configurations for a metal



LLNL Sierra  
(#2 on TOP500 list)



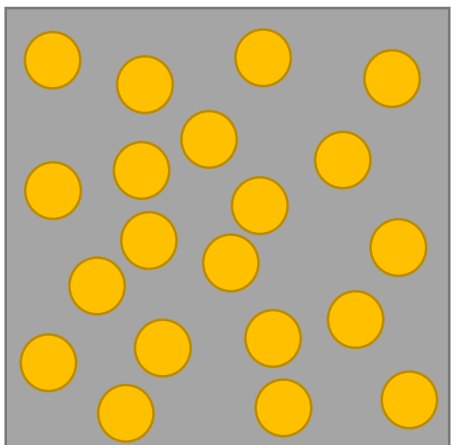
Thanks for the open  
science early access  
allocation!



# Sampling techniques

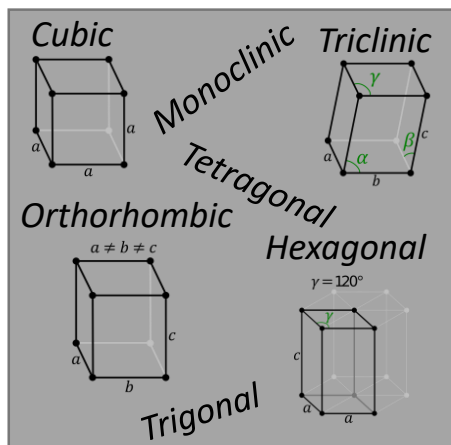
## Random sampling

### Technique 1 Disorder



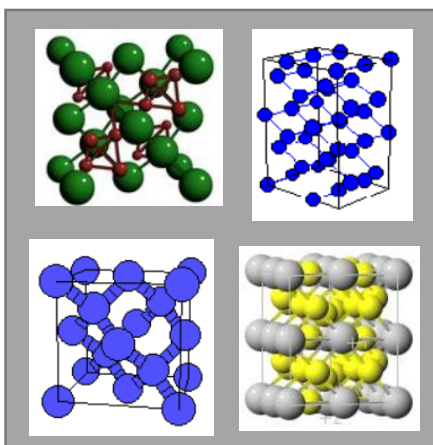
- Configuration selected by ML
- Minimum atomic distances restrained
- Density kept within a set range
- Random a, b, c lattice constant
- Box size kept minimal

### Technique 2 Space group<sup>1</sup>



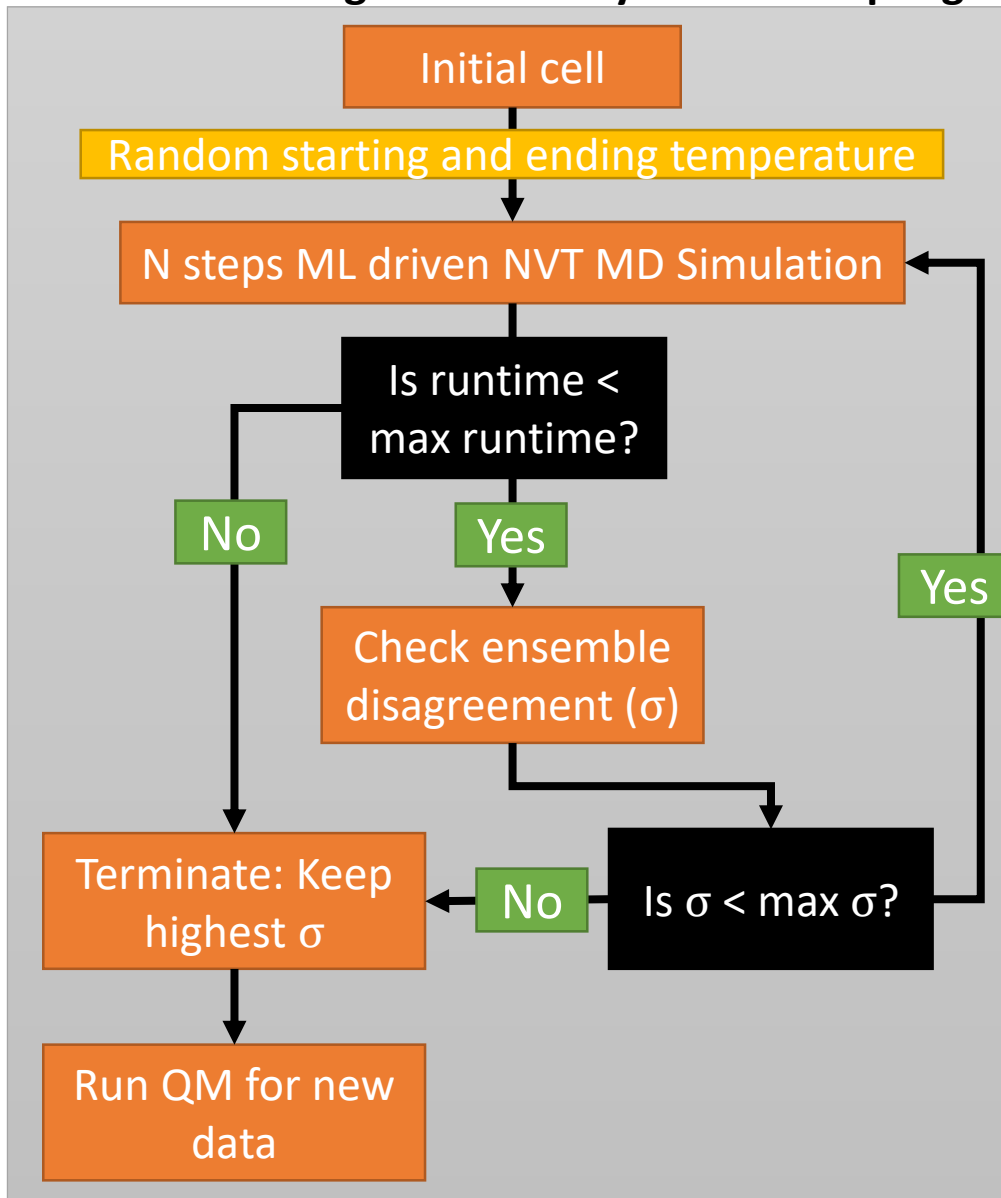
## Selected sampling

### Technique 3 Crystal<sup>2</sup>



- Crystal selected by human
- Random perturbation by 0.25Å
- Random a, b, c lattice constant

## Active learning molecular dynamics sampling



1) Images from: [https://en.wikipedia.org/wiki/Space\\_group](https://en.wikipedia.org/wiki/Space_group)

2) Images from: <https://homepage.univie.ac.at/michael.leitner/lattice/index.html>

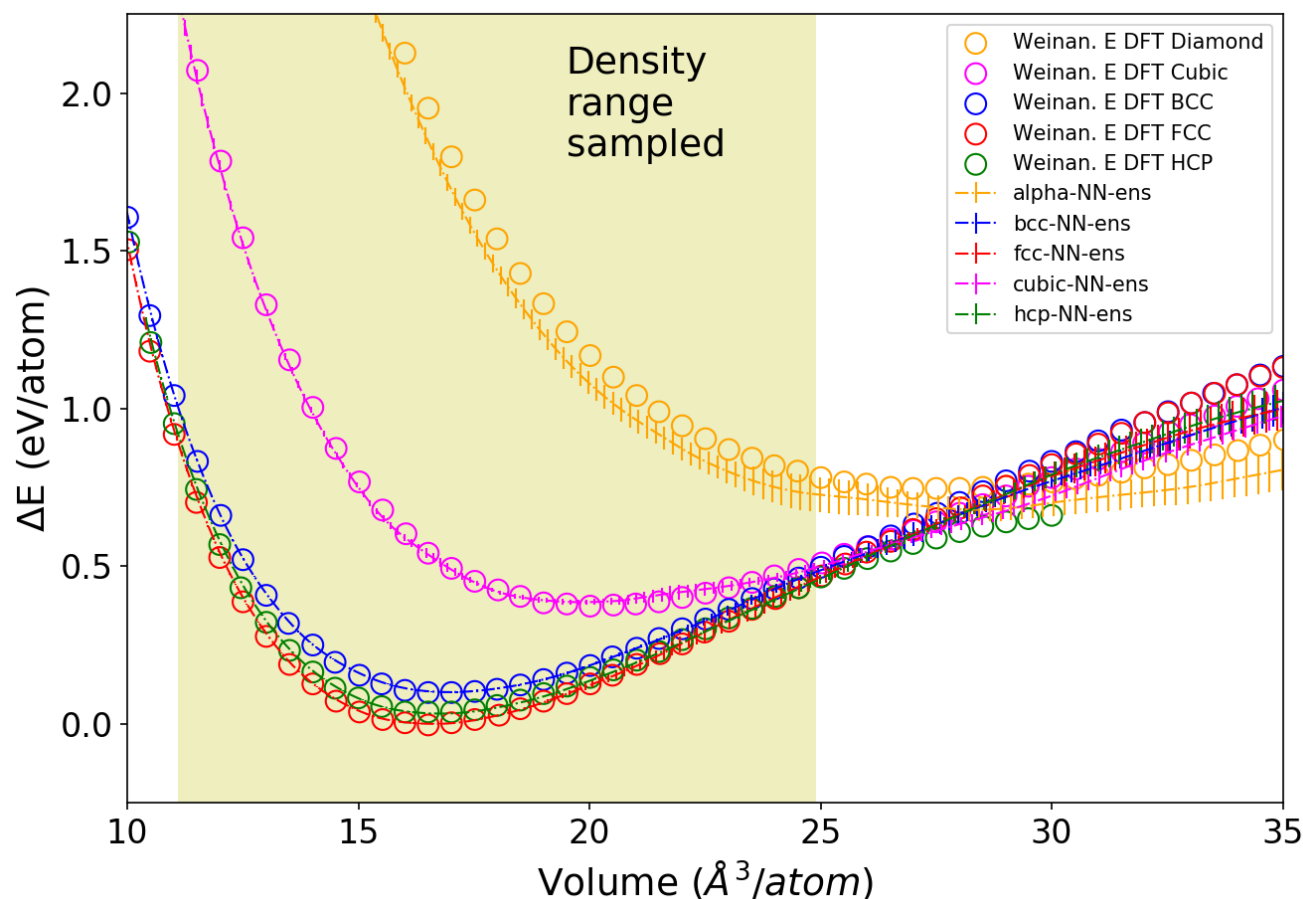
# Application on Aluminum (Al)!

## DISCLAIMER!

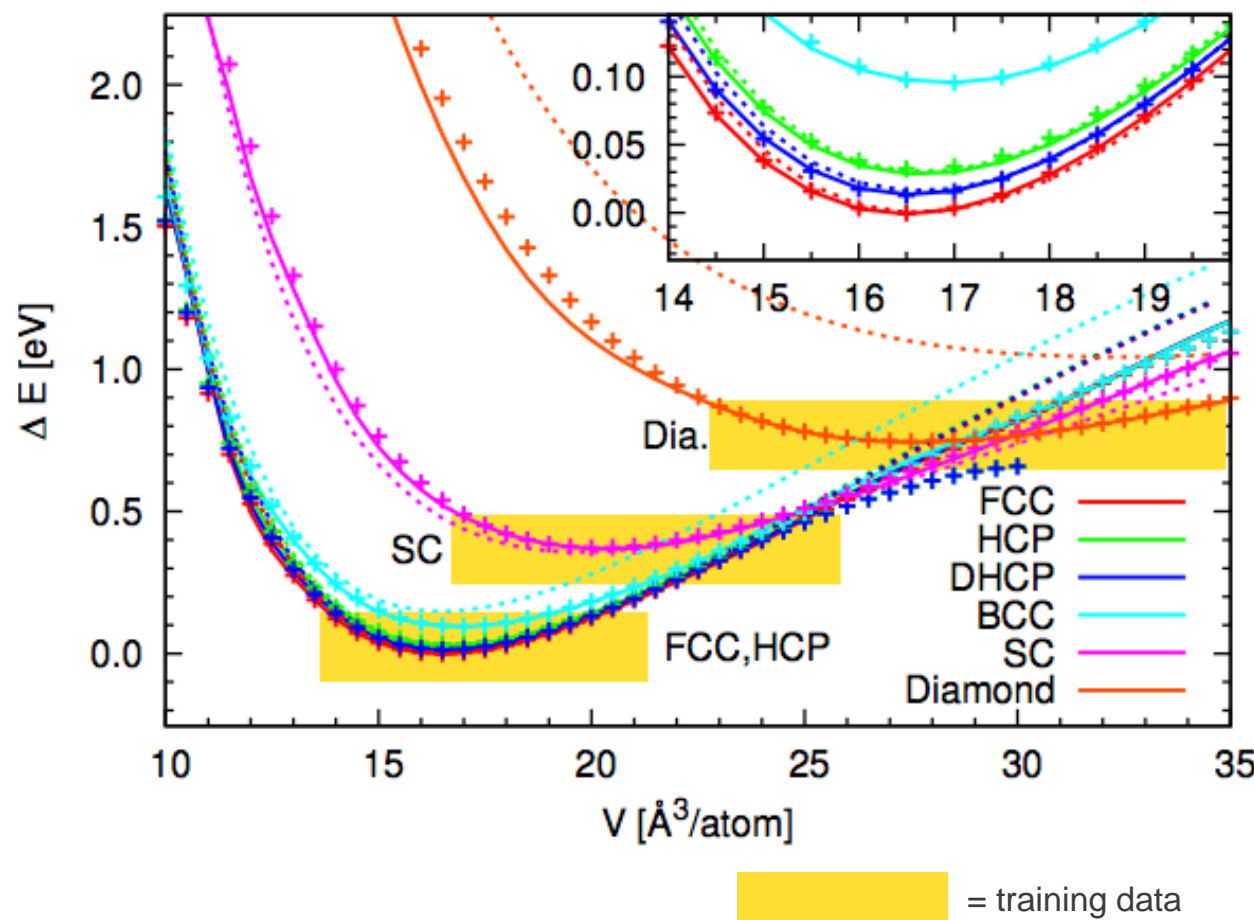
We trained to DFT and are about to show results comparing to experiment.

# Select crystal vs. random disorder MD sampling for AI

**No human knowledge used in sampling**  
**Sampling technique 1 only: Disorder**  
(our current work)



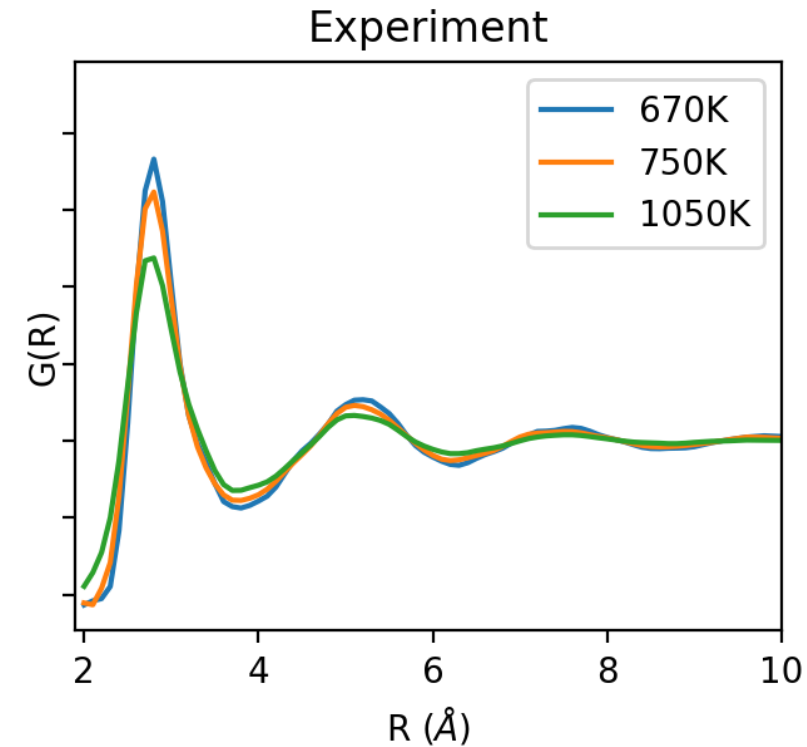
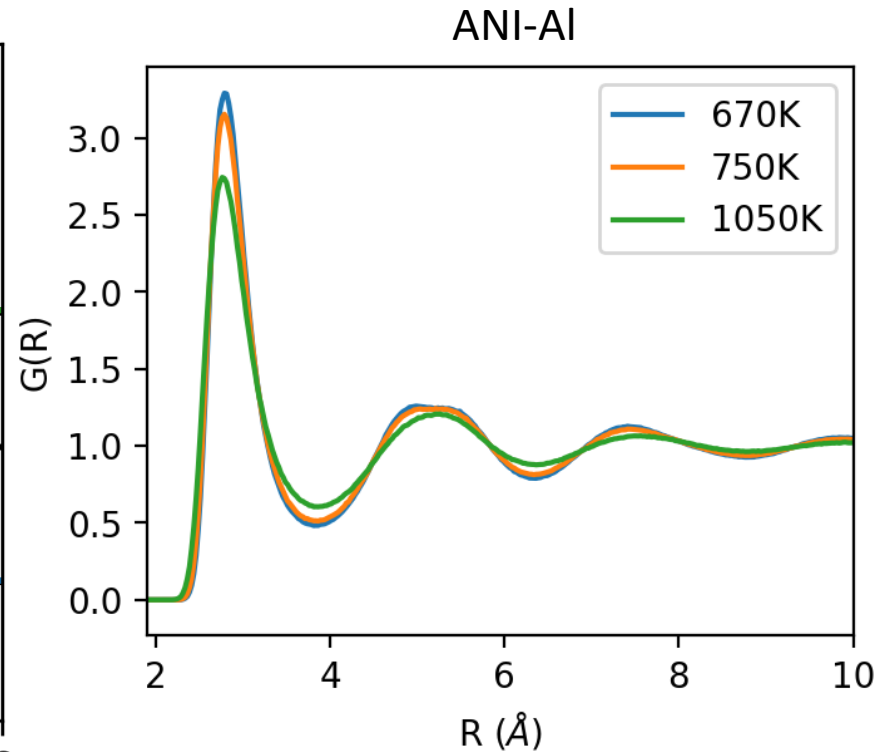
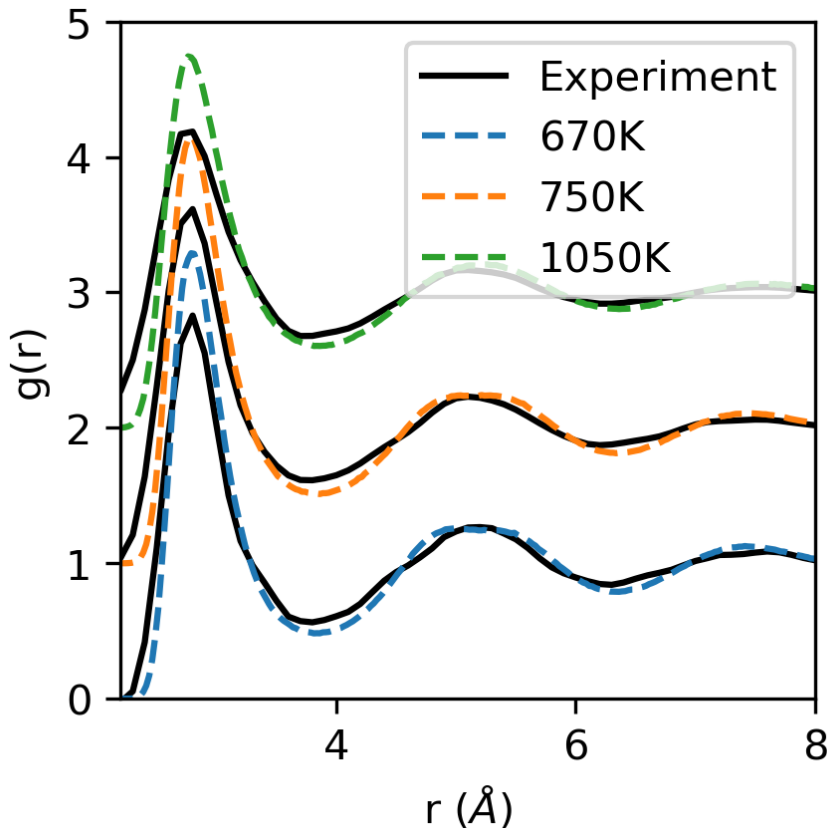
**Crystals chosen based on human knowledge (previous literature)**



Linfeng Zhang, De-Ye Lin, Han Wang, Roberto Car, and Weinan E,  
Active Learning of Uniformly Accurate Inter-atomic Potentials for  
Materials Simulation, [arXiv:1810.11890]

# RDF of liquid Al using our ML potential

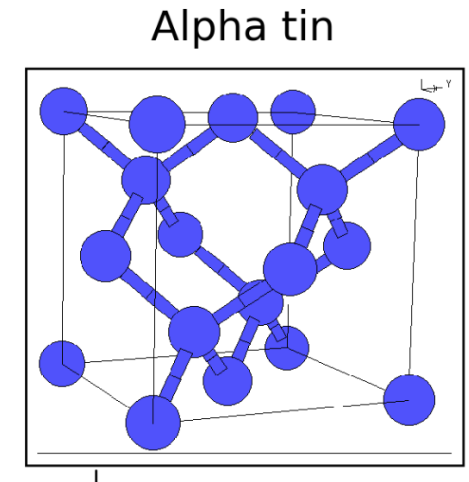
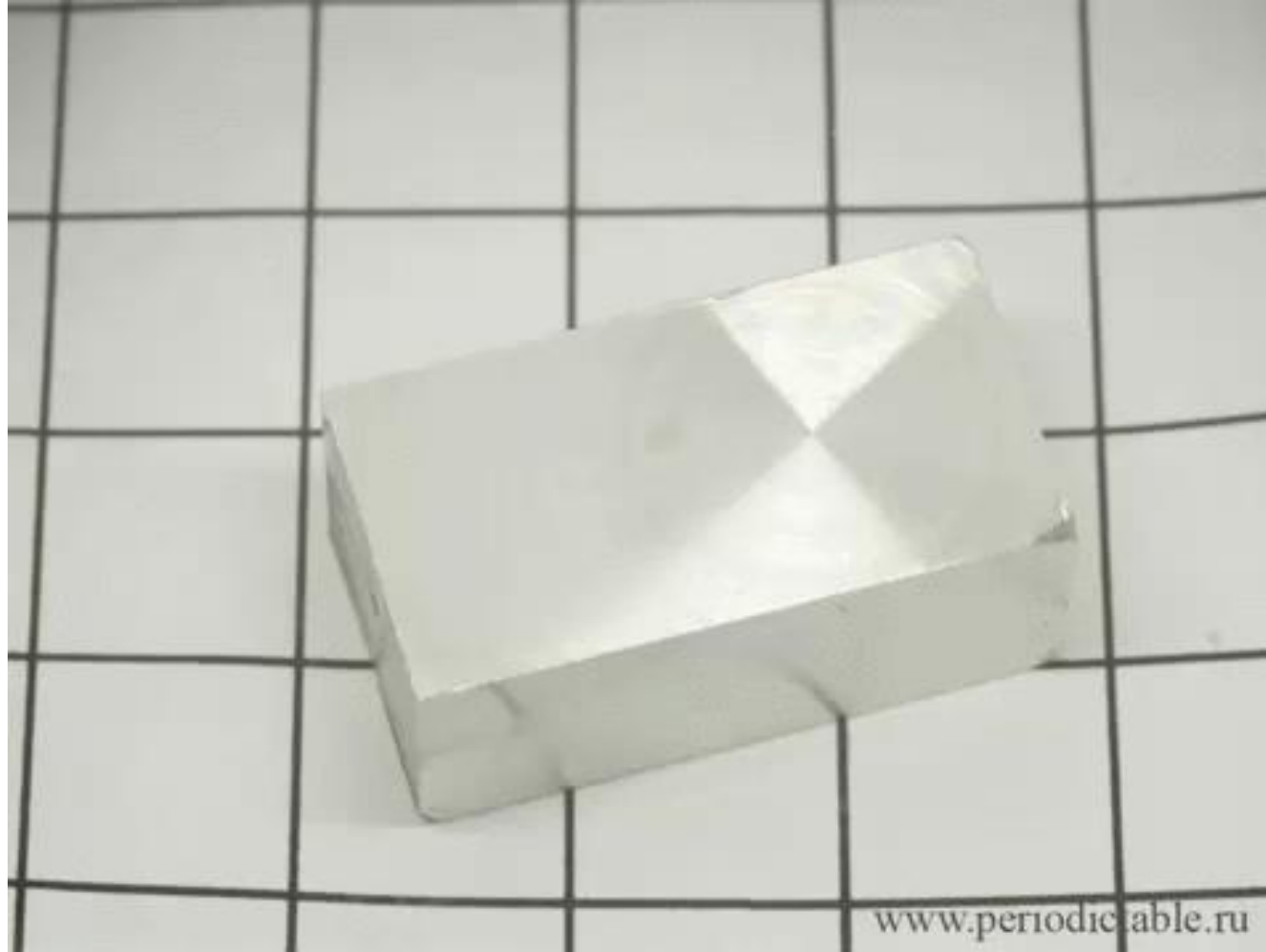
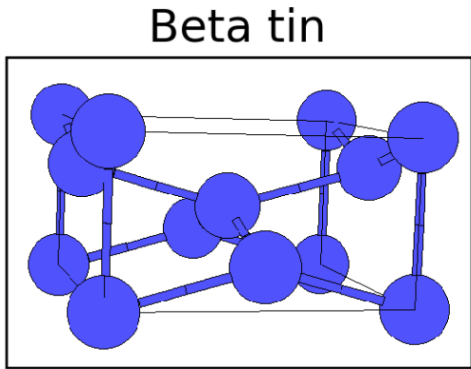
- 125ps of NPT for equilibration
- 125ps of NVT at equilibrated density
- 2048 atoms system
- Trained to DFT (PBE)
- Density vs exp ~11% off (on the level of typical DFT error)



Exp Data: [http://res.tagen.tohoku.ac.jp/~waseda/scm/LIQ/periodic\\_table.html](http://res.tagen.tohoku.ac.jp/~waseda/scm/LIQ/periodic_table.html)

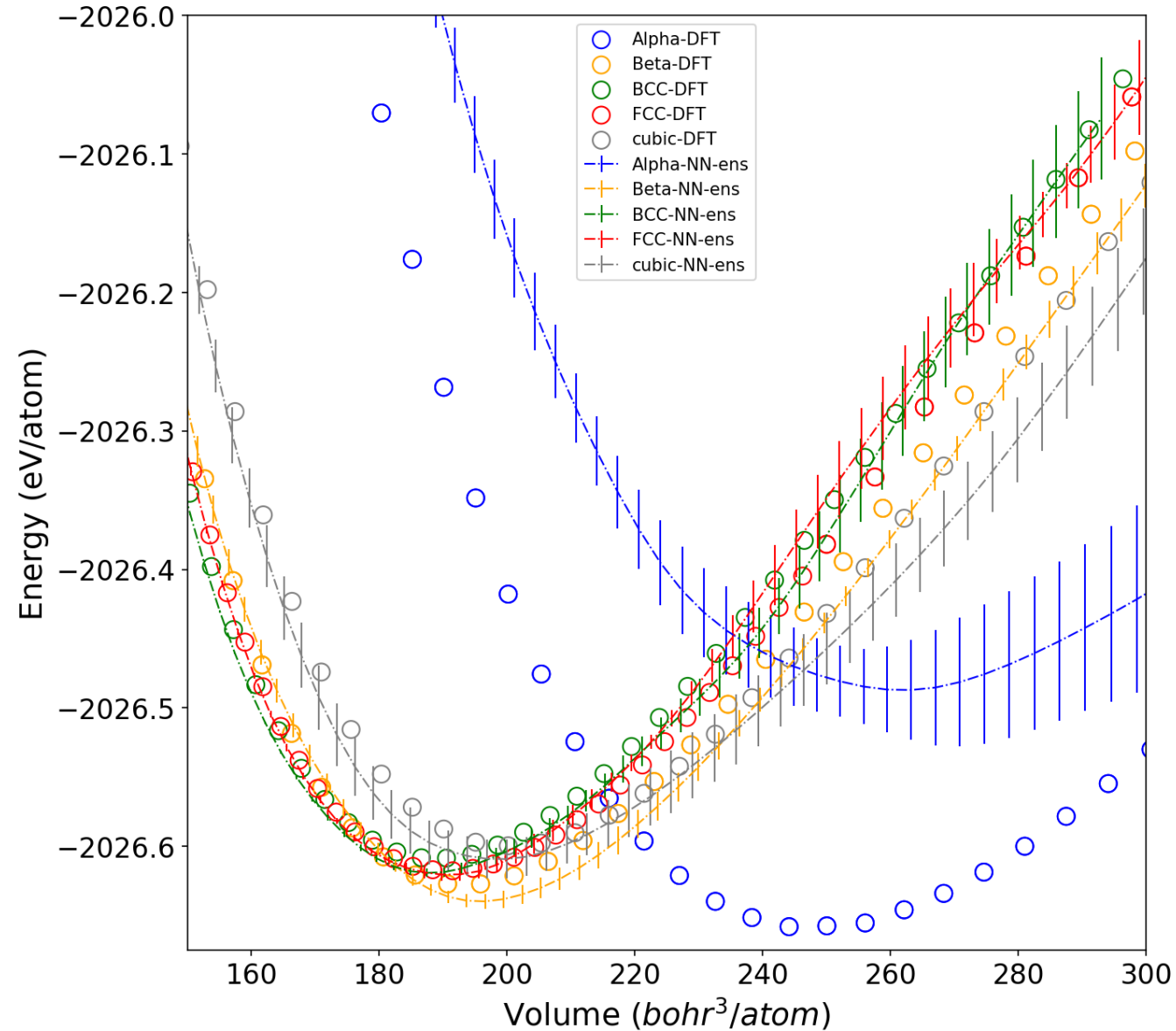
# Application on Tin (Sn)!

## $\beta$ -Sn to $\alpha$ -Sn phase transition

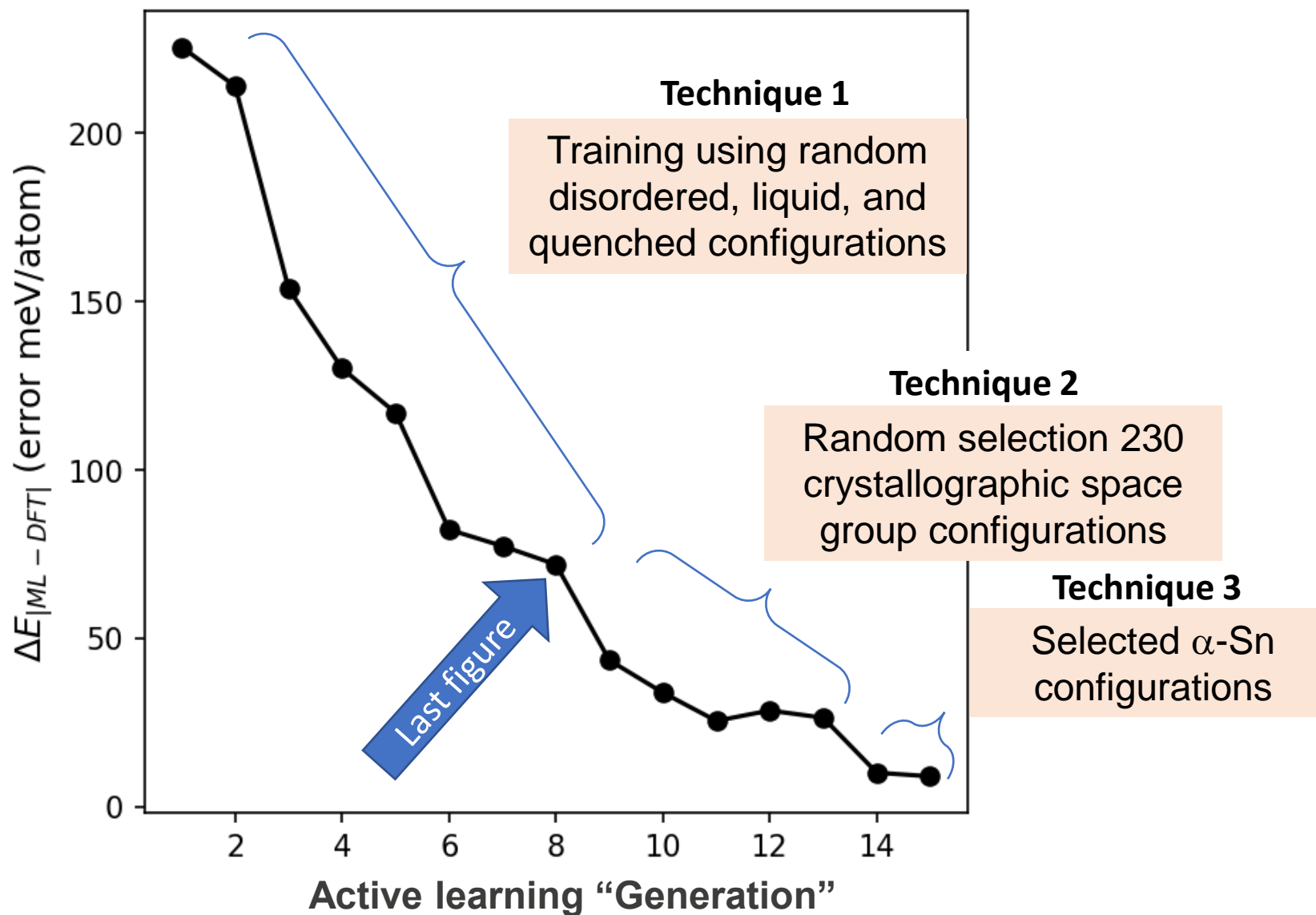


# Random disorder MD sampling for Sn

Active learning on Sn w/o any human intervention



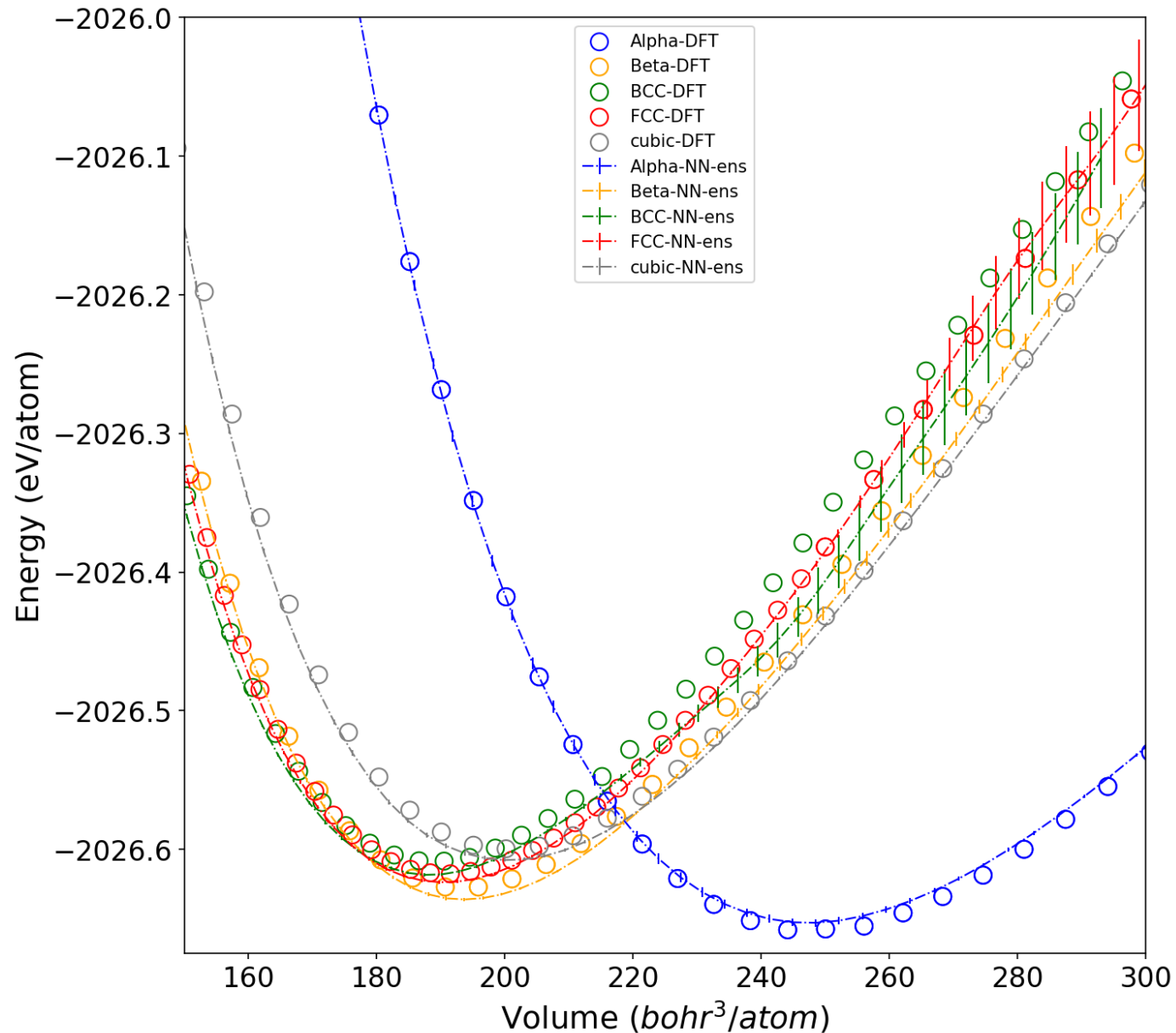
# Error on $\alpha$ -Sn with AL progress



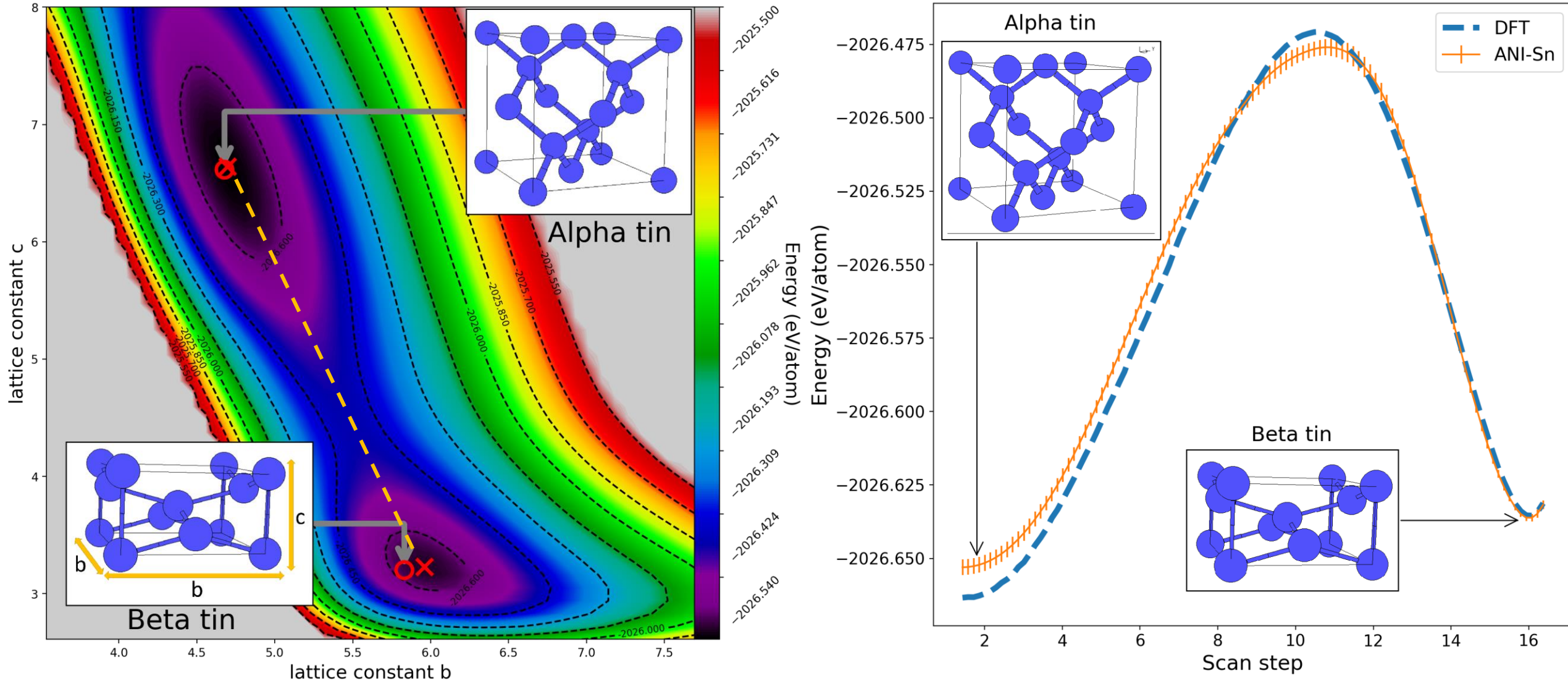


# Alpha Sn requires explicit alpha fitting

Including hand picked crystal structures

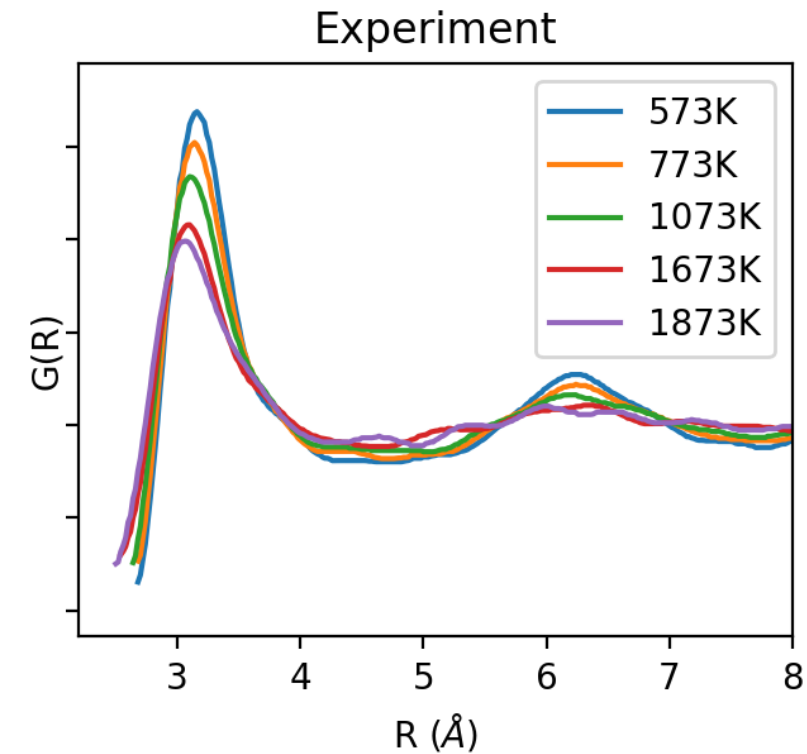
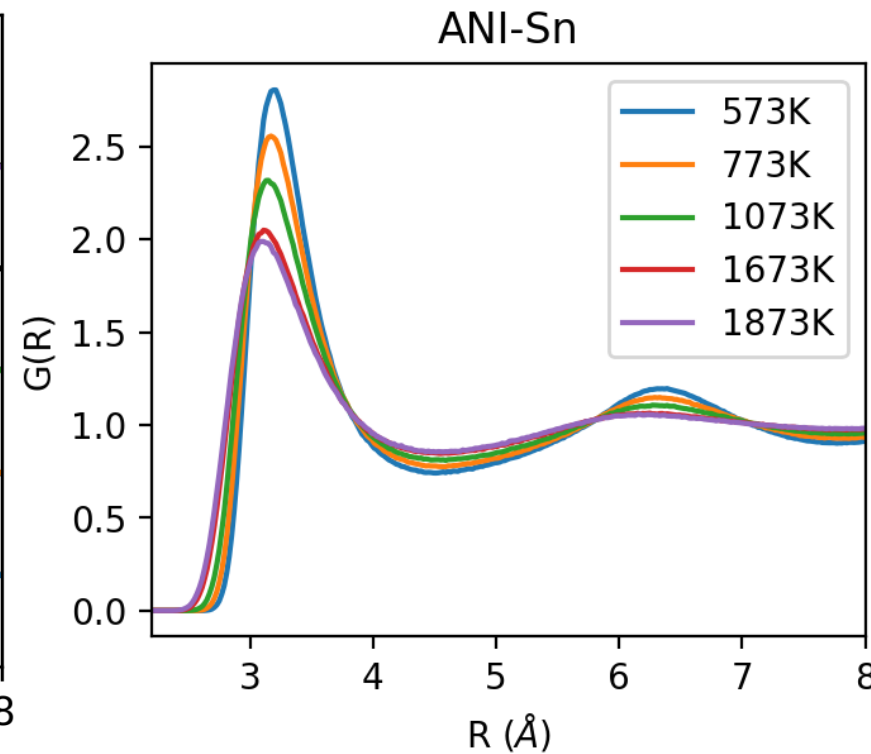
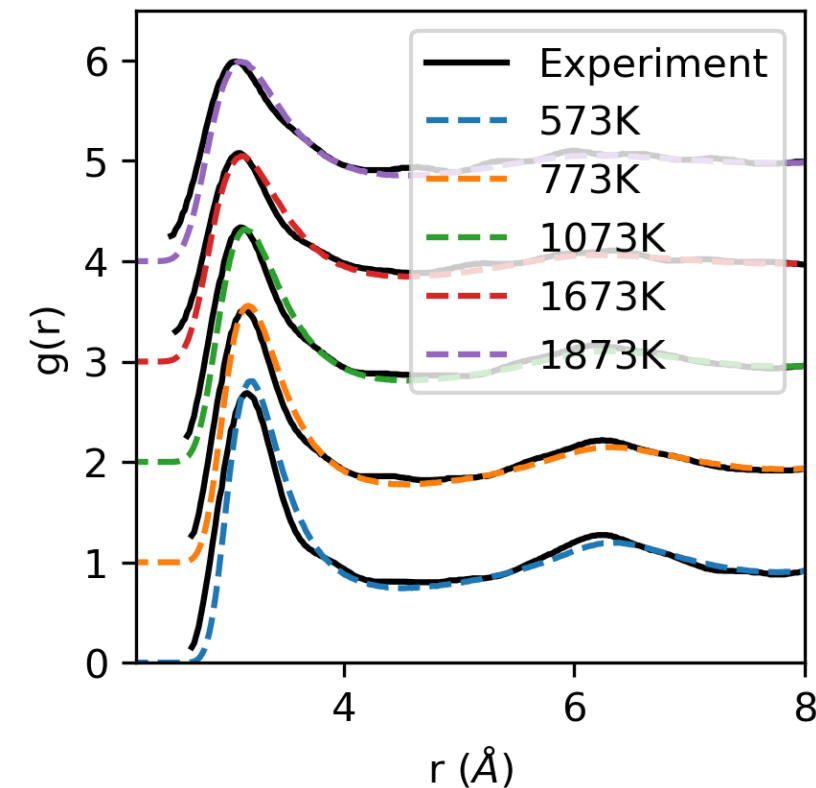


# Good agreement with DFT for crystals AND barriers



# Liquid Sn RDFs at variable temperatures

- 125ps of NPT for equilibration
- 125ps of NVT at equilibrated density
- 1728 atoms system
- Trained to DFT (PBE)
- Density vs exp ~9% off (on the level of typical DFT error)



Exp data: T. Itami, S. Munejiri, T. Masaki, H. Aoki, Y. Ishii, T. Kamiyama, Y. Senda, F. Shimojo, and K. Hoshino *Phys. Rev. B* 67, 064201 (2003)

# How different sampling methods perform

Datasets:

**Random** = Random configurations and random space groups

**Crystal** = Selected randomly perturbed crystals using human knowledge

Typical MEAM fitness on Crystals

Energy RMSE (meV/atom)	Force RMSE (eV/A)
44.0	0.97

## Testing on Crystals

Training Data	Testing Data	Energy RMSE (meV/atom)	Force RMSE (eV/A)
Crystal	Crystal	3.6	0.04
Random	Crystal	13.7 (17.7/8.3)	0.05 (0.05/0.04)
Random + Crystal	Crystal	4.4	0.03

## Testing on Random

Training Data	Testing Data	Energy RMSE (meV/atom)	Force RMSE (eV/A)
Crystal	Random	250.1	1.88
Random	Random	5.9	0.09
Random + Crystal	Random	5.1	0.09

# Conclusion and Outlook

## Conclusions

- Sampling matters in dynamical studies
- Better data makes a better ML potential
- Active learning methods are required for better data
- Current models may be missing the ability to describe physics in some metals
- Test set results for atomistic ML models can be misleading

## Opportunities

- Continue to develop better sampling techniques for metals and molecules
- Discover better uncertainty quantification methods for active learning
- Apply models to gain physical insights
- Recover long range interactions through combined charge prediction and coulomb models



# Thank you!

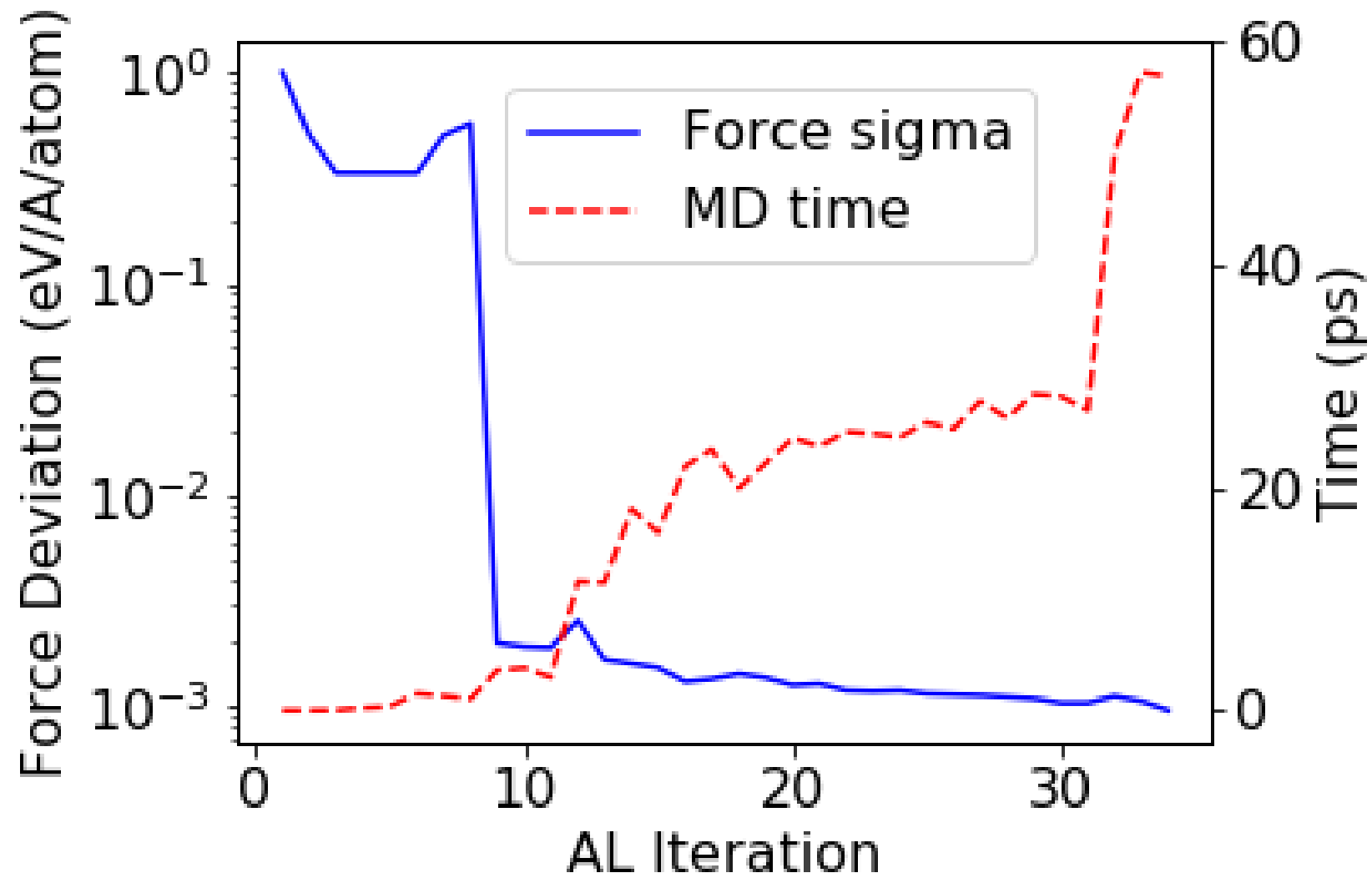


LLNL Sierra  
(#2 on TOP500 list)

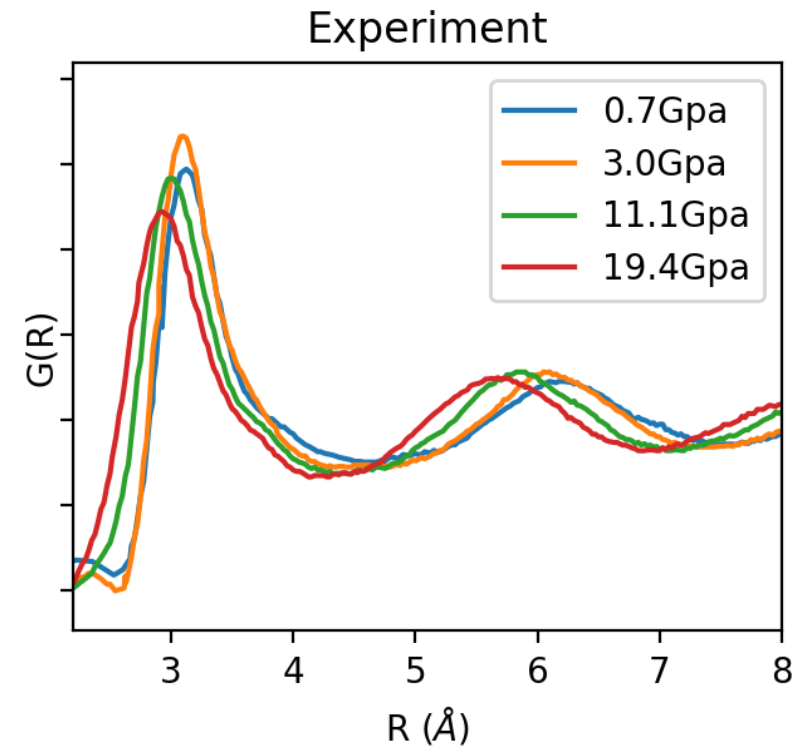
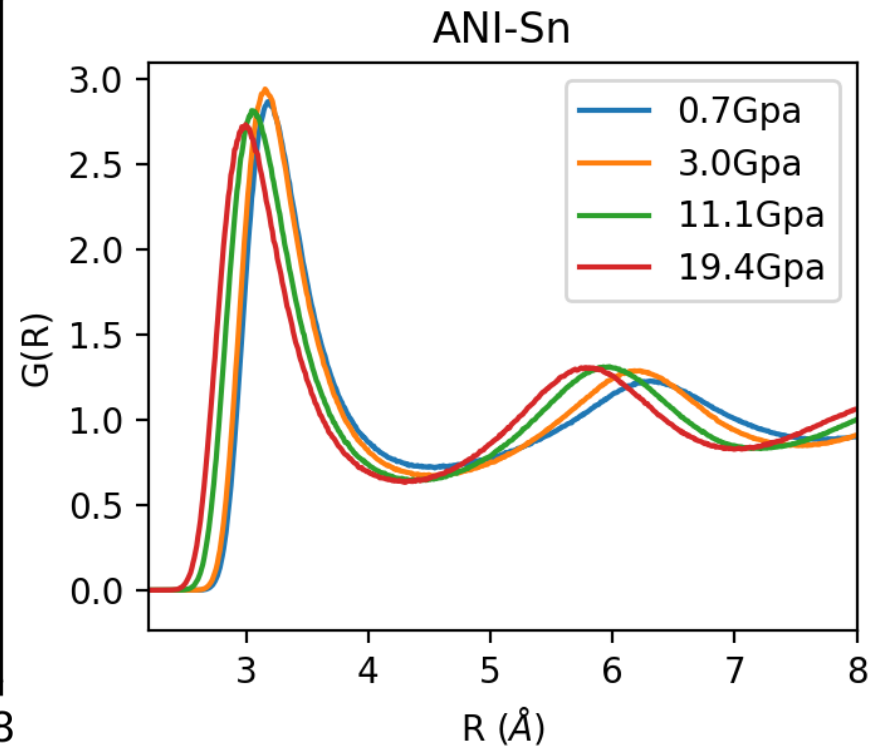
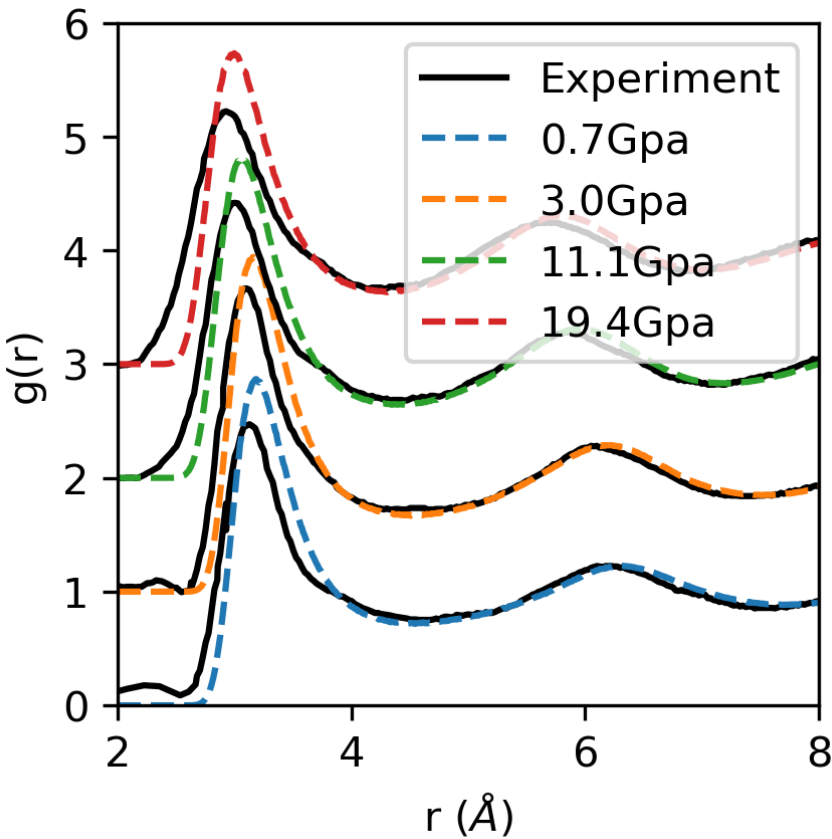


Thanks for the open  
science early access  
allocation!



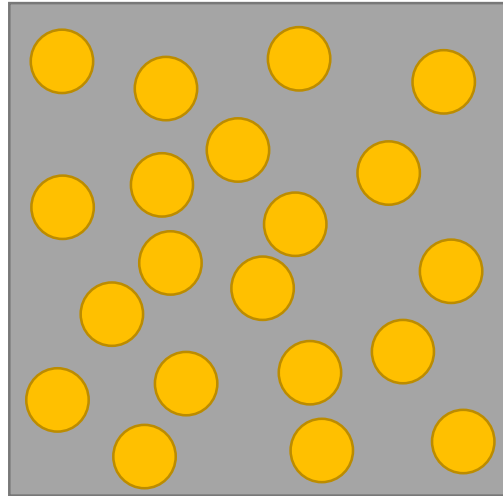


# RDFs at variable pressures



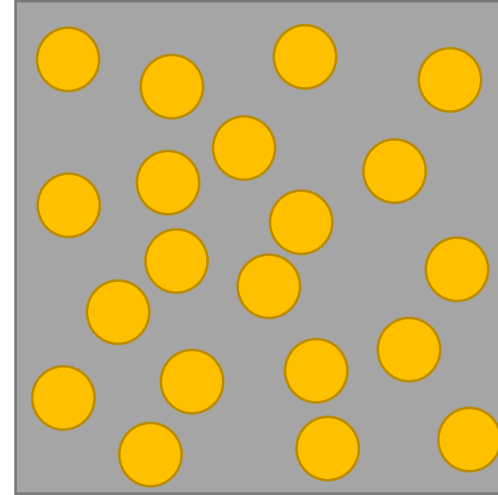
# Random disorder (no human knowledge) sampling technique

Random atom placement



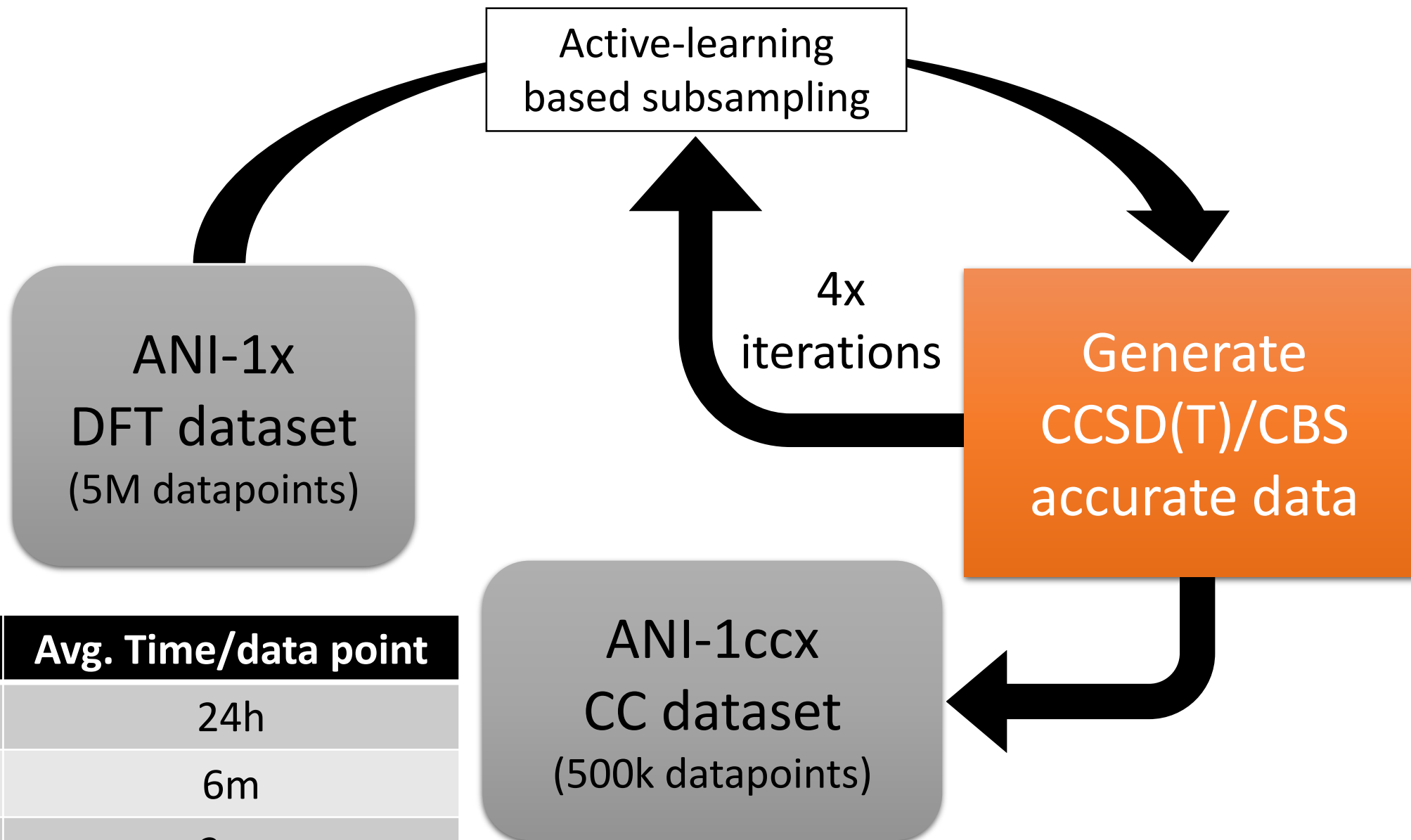
- Minimum atomic distances restrained
- Density kept within a set range
- Random a, b, c lattice constant
- Box size kept minimal

AL MD Sampling with current ML potential



- NVT dynamics
- Randomized starting and ending temperature
- Simulation ends with high ensemble disagreement

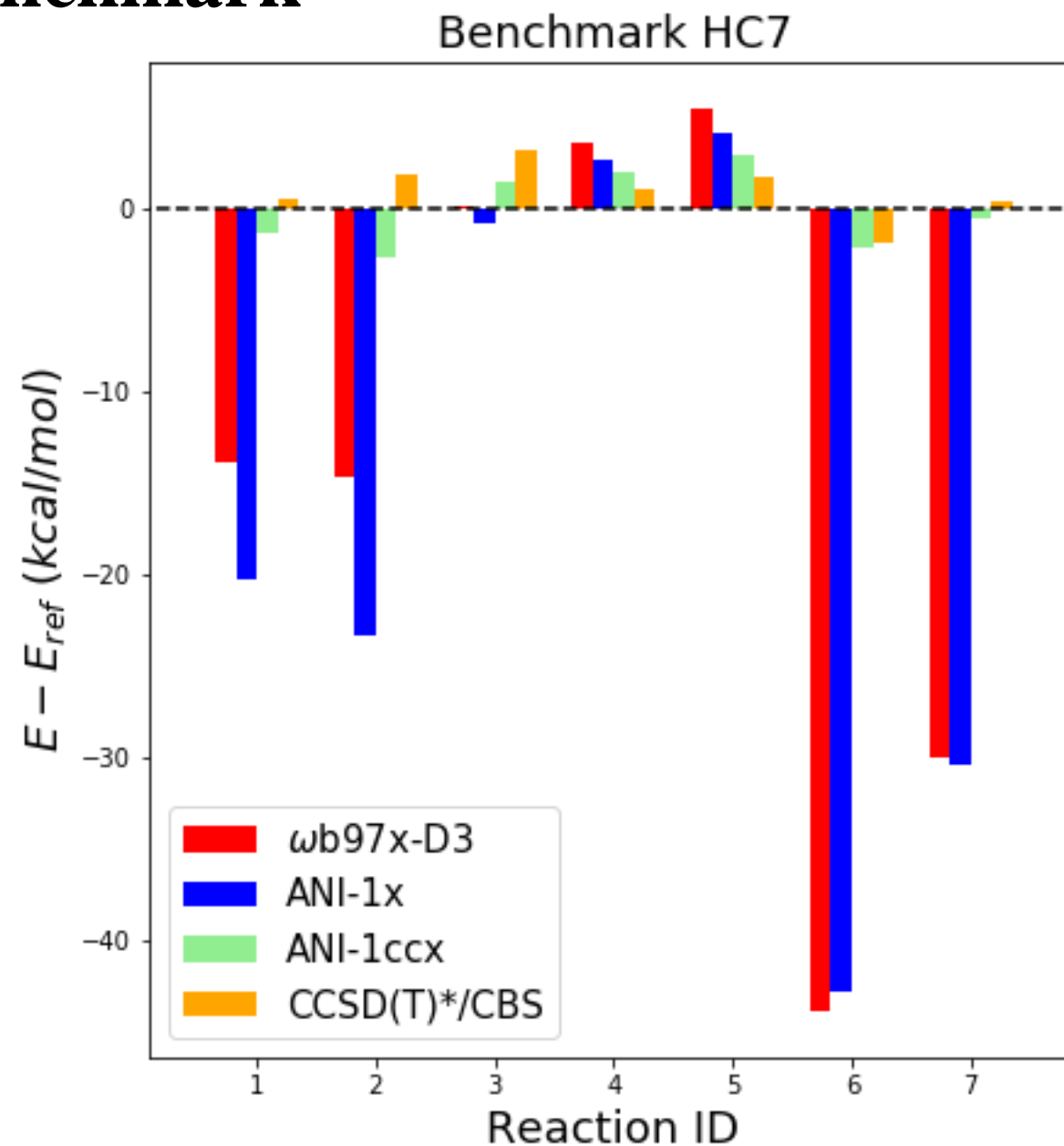
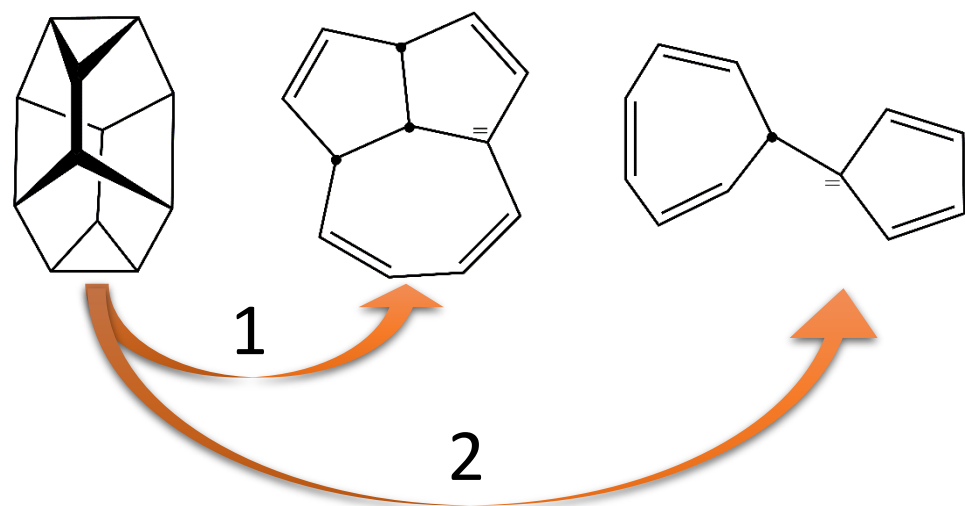
# CCSD(T)/CBS accurate data generation



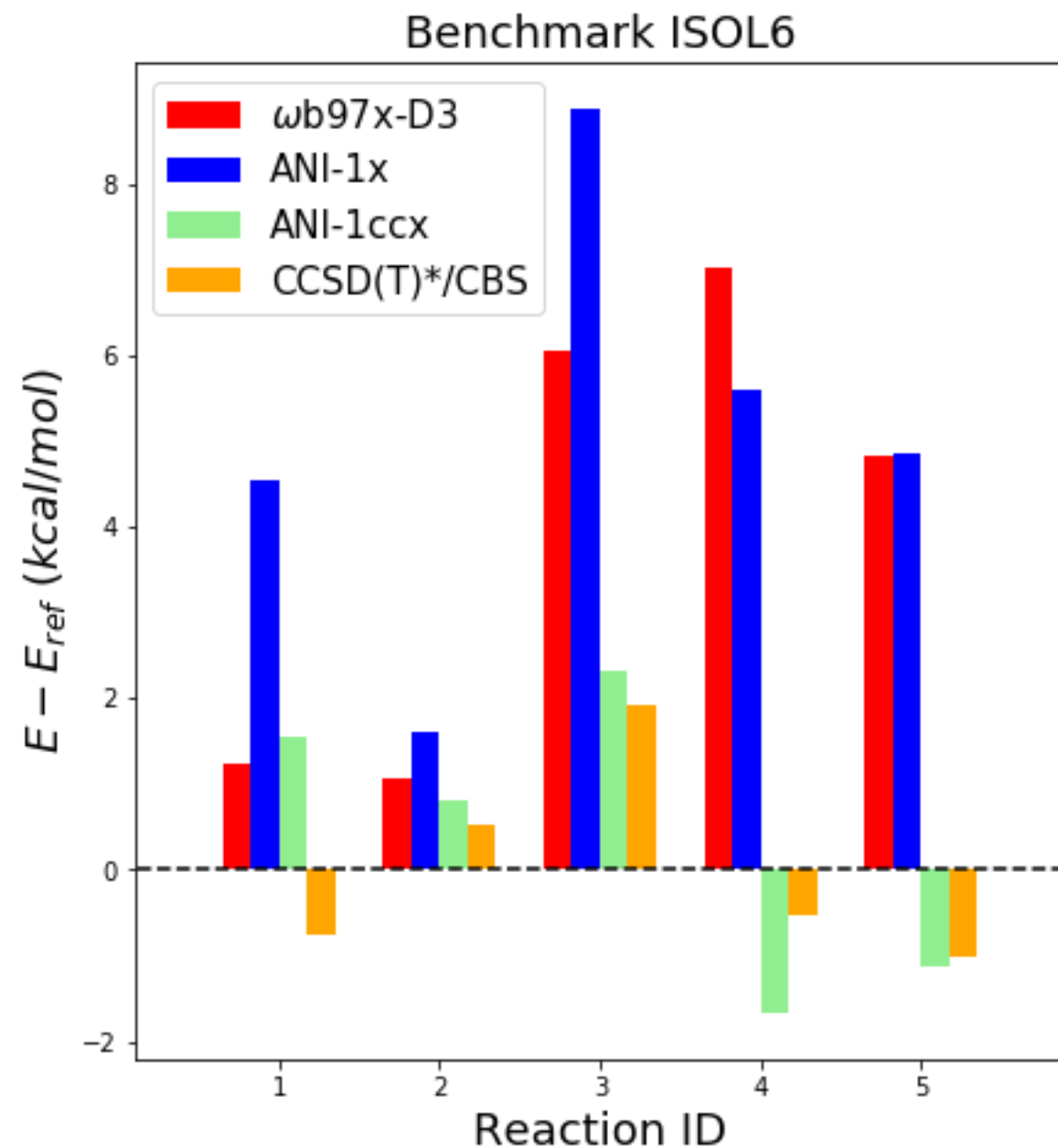
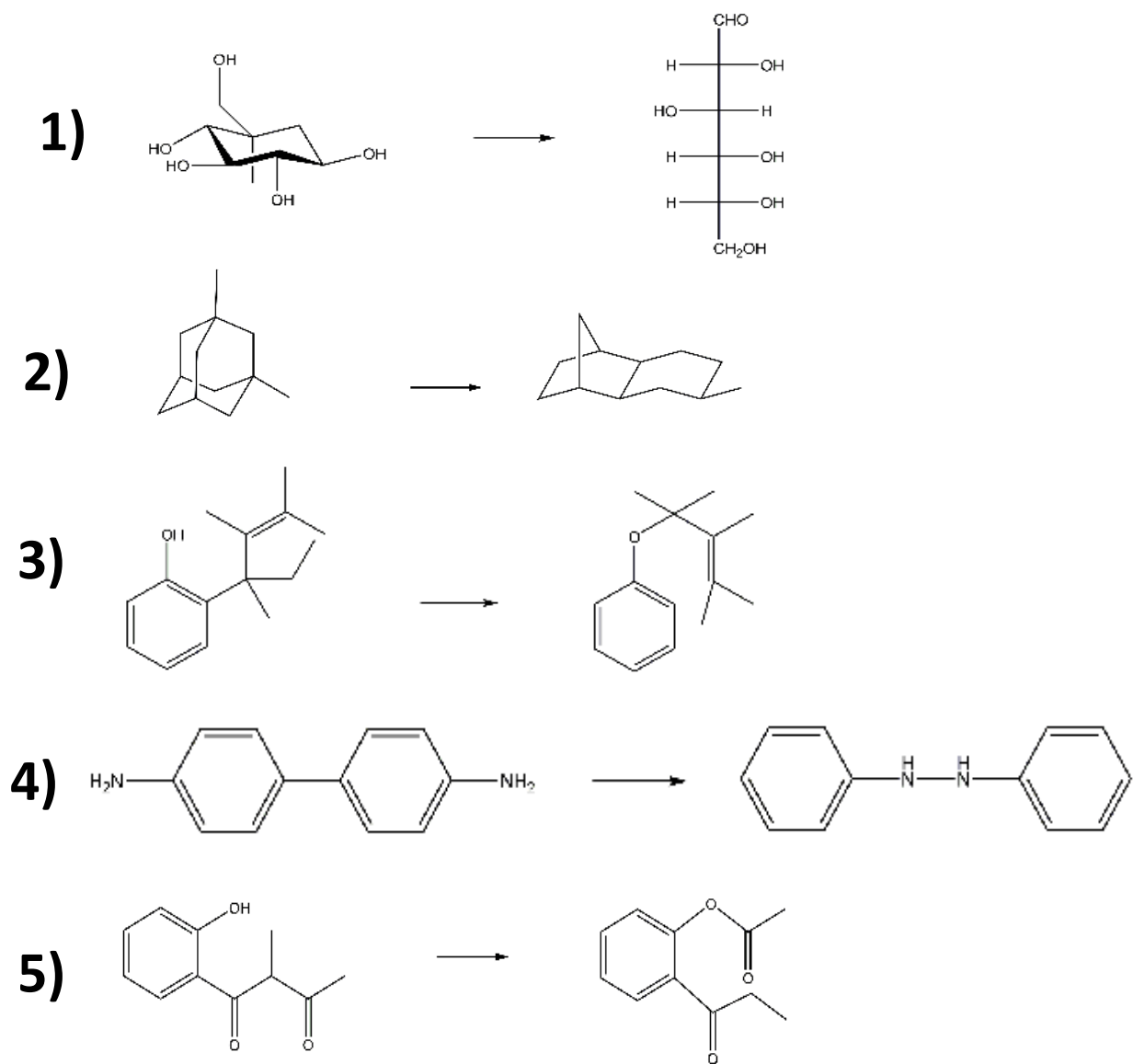
Method	Avg. Time/data point
CCSD(T)	24h
DFT	6m
ANI-1ccx	2μs

# Hydrocarbon reaction energy benchmark

## Examples

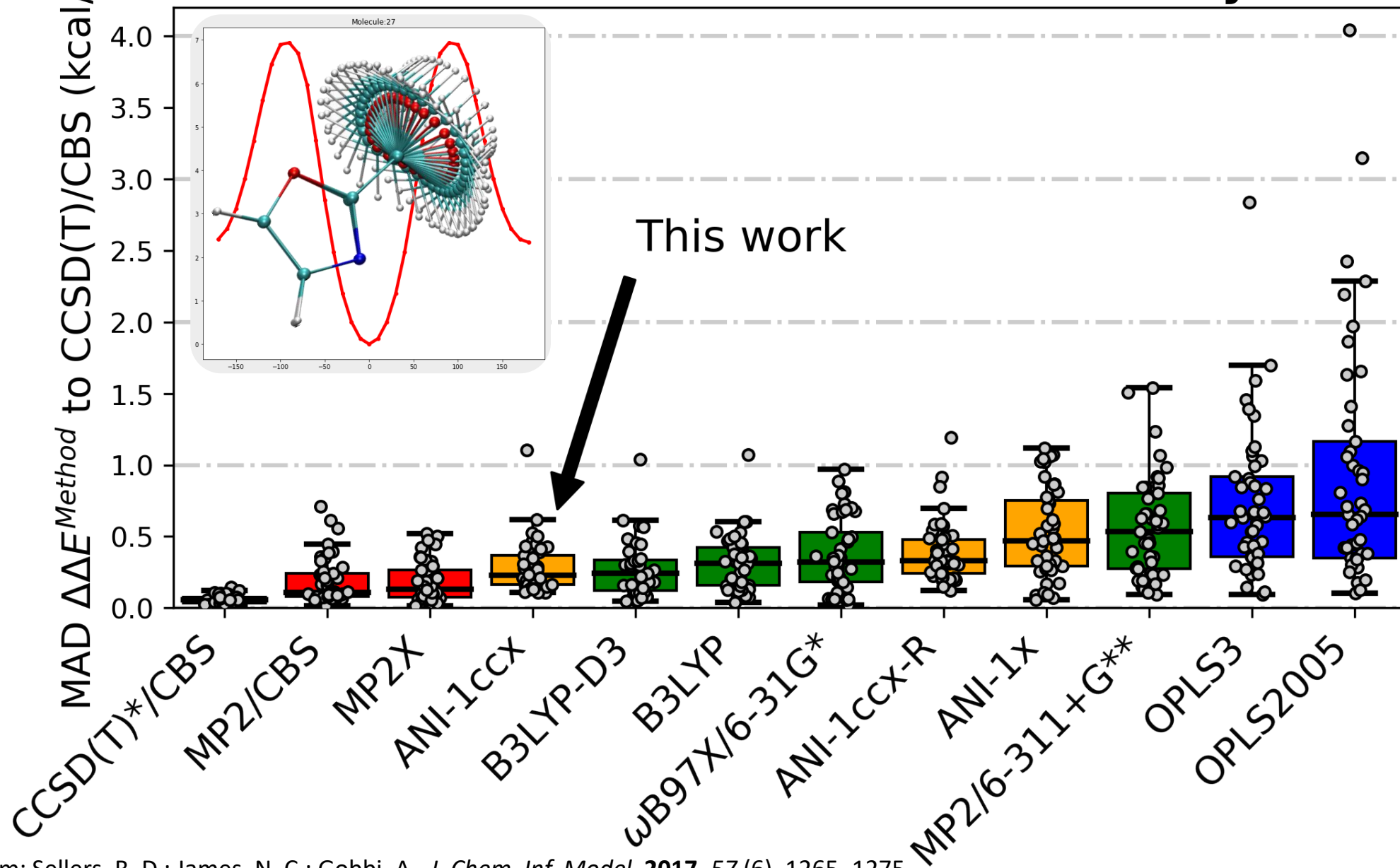


# Organic reaction energy benchmark

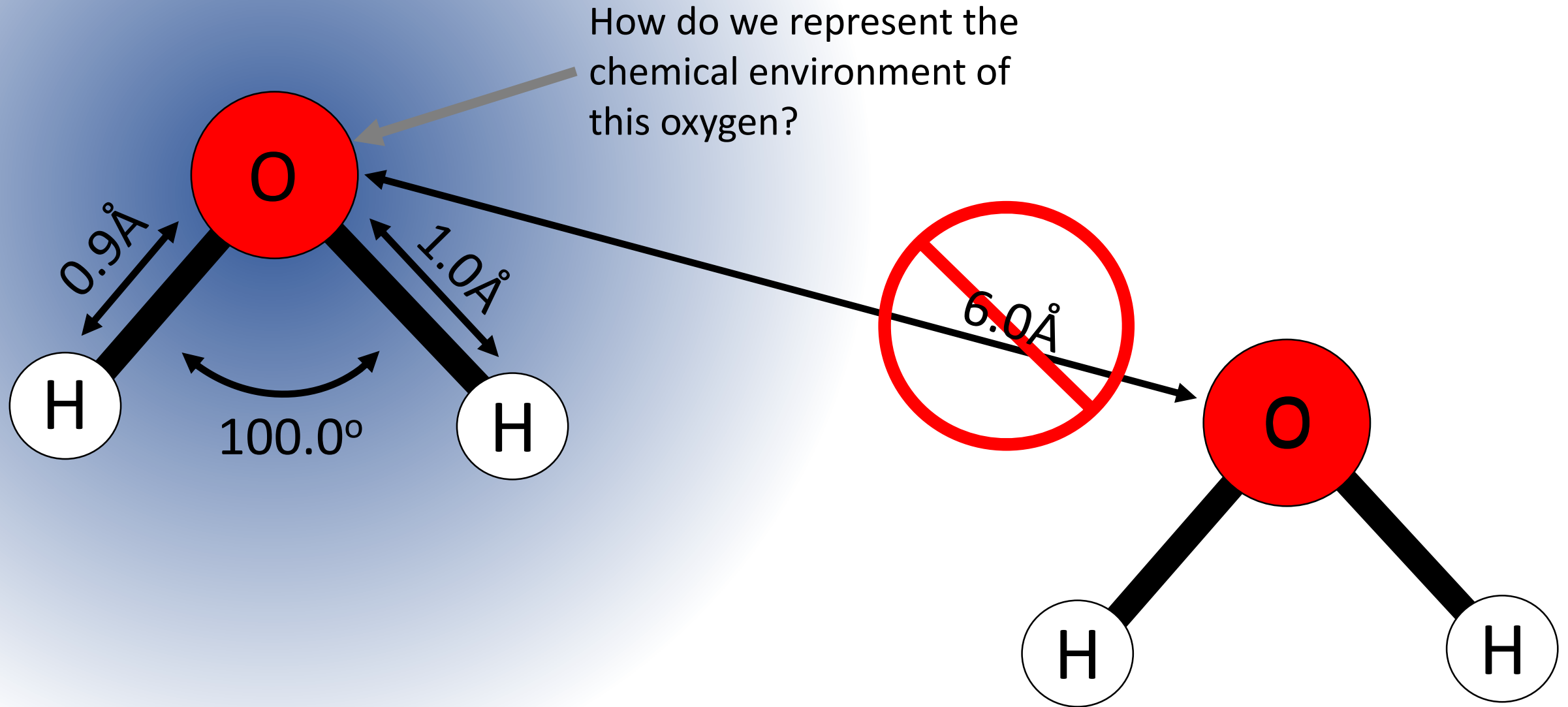




# Torsion benchmark (CHNO only)



# Atomic environment description



# Descriptors for the ANI ML-based potential

## Radial symmetry function

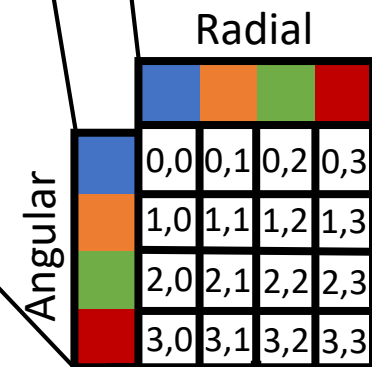
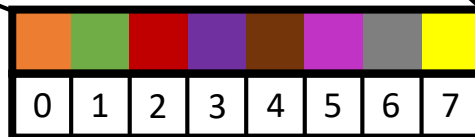
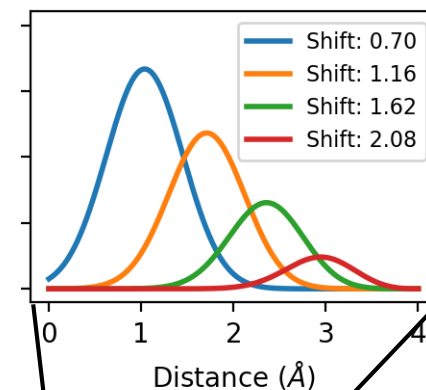
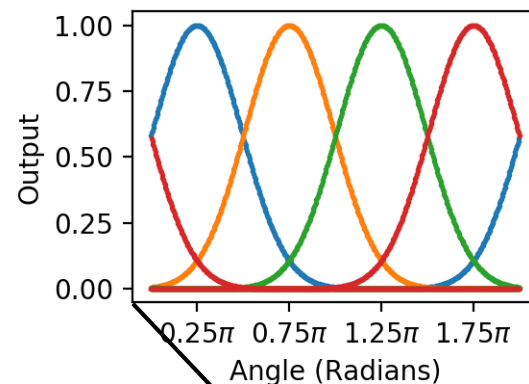
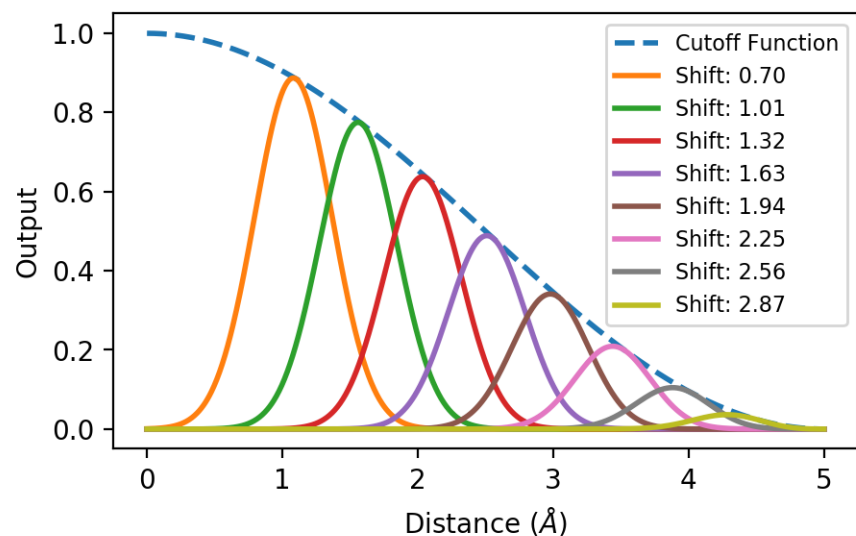
$$G_m^R = \sum_{j \neq i}^{\text{All Atoms}} e^{-\eta(R_{ij}-R_s)^2} f_c(R_{ij})$$

## Cutoff function

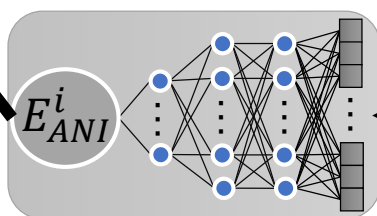
$$f_c(R_{ij}) = \begin{cases} 0.5 \times \cos\left(\frac{\pi R_{ij}}{R_c}\right) + 0.5 & \text{for } R_{ij} \leq R_c \\ 0.0 & \text{for } R_{ij} > R_c \end{cases}$$

## Angular symmetry function

$$G_m^{A_{\text{mod}}} = 2^{1-\zeta} \sum_{j,k \neq i}^{\text{All Atoms}} (1 + \cos(\theta_{ijk} - \theta_s))^\zeta \exp\left[-\eta\left(\frac{R_{ij}^2 + R_{ik}^2}{2} - R_s\right)^2\right] f_c(R_{ij}) f_c(R_{ik})$$



$$E = \sum_i^N E_i$$



Concatenate